

# Leveraging Large Language Models for In-context Data Generation to Address Bias and Scarce Data Challenges in Mental Healthcare

Anonymous ACL submission

## Abstract

Motivational Interviewing (MI) is a widely used, evidence-based counseling approach, yet the development of robust NLP systems for MI remains constrained by severe data scarcity, high annotation costs, and the need for domain expertise. While recent work has explored large language models (LLMs) for clinical text generation and augmentation, existing studies largely focus on utterance-level transformations or evaluate a limited set of models, leaving the role of long-range conversational context underexplored. In this work, we present the first systematic, session-level benchmark of LLM-based data augmentation for MI dialogues. We compare prompt-based generative augmentation at the full-session level with utterance-level task-sensitive transformation methods across 13 state-of-the-art (SOTA) LLMs and three long-context classification models designed for extended clinical conversations. As a key outcome, we also present **ICAUGAnnoMI**, a novel dataset of 1,764 low- and high-quality MI dialogues, spanning nearly 81k talk turns between therapist and client alongside a fidelity-aware evaluation framework that assesses semantic drift, hallucination, and adherence to core MI principles. Our empirical results demonstrate that session-level augmentation consistently outperforms utterance-level approaches in improving MI session classification, particularly for minority and low-quality cases, highlighting the importance of preserving long-range conversational structure. Beyond performance gains, our analysis provides practical insights into the strengths and limitations of contemporary LLMs in sensitive mental health settings. We release our code and commit to publicly sharing the generated data to facilitate reproducible research and future benchmarking in MI-focused NLP. ICAUGAnnoMI data and source code is available at **repository**<sup>1</sup>.

<sup>1</sup>[https://anonymous.4open.science/r/ARR\\_Submission\\_Cycle\\_2026/README.md](https://anonymous.4open.science/r/ARR_Submission_Cycle_2026/README.md)

## 1 Introduction

Anxiety and depression are among the most common psychological disorders, affecting social and economic functioning. According to the latest WHO report<sup>2</sup>, mental disorders account for 1 in 6 years lived with disability, and individuals with severe conditions die 10–20 years earlier than the general population (Organization, 2022). AI-driven solutions have recently gained attention for alleviating this burden by supporting early screening, psycho-education, and self-help strategies (Shiffrin and Mitchell, 2023; Stella et al., 2023; Demszky et al., 2023; Sezgin, 2023; Balloccu et al., 2024a). Large Language Models (LLMs) are particularly promising due to their human-like dialogue capabilities (Han et al., 2024; Kumar et al., 2024, 2025) and potential to improve access in low-resource settings (Xu et al., 2024). However, integrating LLMs into mental healthcare faces challenges (Wang et al., 2023): i) *hallucination*, generating plausible but incorrect text (Li et al., 2023a); ii) *stochastic parroting*, producing fluent but shallow responses (Bender et al., 2021); iii) *bias*, reducing reliability for underrepresented groups due to representation and sampling issues (Balloccu et al., 2024b; Badyal et al., 2023; Morales et al., 2023; Harrer, 2023). These challenges converge on a single issue: the *lack of publicly available high-quality data for effective training*. To address this, we employ state-of-the-art (SOTA) LLMs for in-context human-like data generation at the MI session level (Brown et al., 2020; Chen et al., 2022; Dong et al., 2022) and task-sensitive NLP augmenters at the utterance level (Dhole et al., 2023).

Our experimental study is based on AnnoMI (Wu et al., 2022, 2023; Kumar et al., 2023), a publicly available expert-annotated dataset com-

<sup>2</sup><https://www.who.int/teams/mental-health-and-substance-use/world-mental-health-report>



ones. We define a feature mapping  $\Gamma : \mathcal{S} \rightarrow \mathbb{R}^d$ , a classifier  $f_\theta : \mathbb{R}^d \rightarrow [0, 1]$ , and the following partition  $\mathcal{D} = \mathcal{D}_{train} \cup \mathcal{D}_{test}$ . Let  $\phi$  be an augmentation function, drawn from some family of augmentation functions  $\Phi$ . Formally:

$$\phi : \mathcal{D}_{train} \mapsto \mathcal{D}_{aug} \quad (1)$$

that is  $\phi$  generates an augmented variation  $\mathcal{D}_{aug}$  of the original training set. We seek to find the augments  $\phi^*$  that solves the empirical risk minimization problem as shown in Appendix C.

## 2.2 Dataset

For this study, we used the GDPR-compliant dataset **AnnoMI**, which consists of 131 low- and high-quality MI dialogues, each containing talk turns between the client and the therapist (Wu et al., 2022, 2023; Kumar et al., 2023), for a total of 9,649 utterances. The dataset contains 23 low-quality and 108 high-quality MI dialogues. AnnoMI is among the few publicly available datasets that contain labeled MI sessions at both the talk-turn and entire conversational levels. This dataset orientation complements our experimental setup of synthetic MI generation and, thus, is the obvious choice for this study. For ease of understanding, Figure 1 in Appendix B depicts a sample of low- and high-quality MI sessions.

## 2.3 Data augmentation

We specifically focus on generating human-like and contextually sound MI sessions that maintain therapeutic fidelity and reflect the diverse range of patient responses and therapist strategies. Given the high cost and complexity of acquiring large-scale MI data, we explore two complementary strategies for data augmentation. First, we employ task-sensitive transformations, which alter transcripts through techniques specialized in text manipulation, e.g., paraphrasing. Second, we utilize prompt-based generation, leveraging LLMs to create new, contextually accurate MI dialogues.

### 2.3.1 Task-sensitive transformation

Task-sensitive transformations (TTs) refer to augmentation methods that adopt text-driven modifications on textual data with a specific expected output in mind, such as enhancing linguistic diversity, shifting stylistic properties, or rephrasing. These strategies are typically adopted for analyzing the impact of text manipulations on trained models, both from a negative (concerning text perturbation)

and a positive (concerning text augmentation) viewpoint. Considering that their diffusion precedes the advent of LLMs, they are not typically devised for transforming long input sequences. In our scenario involving long MI sessions, it is counterintuitive to adopt TTs on several client-therapist talk-turns, which would cause the underlying language models (LMs) to drop excessive text and provide trimmed outputs. It follows that TTs should only be adopted at the utterance level to generate augmented data that preserves the original topic the most. Despite the downside of overlooking the contextual information of entire sessions, TTs represent solid augmentation baselines, ranging from rich paraphrasing variety to subtle word edits. Let  $\Lambda : \mathcal{T} \rightarrow \mathbb{R}^e$  be a feature mapping for individual talks  $t \in \mathcal{T}$ , we can formally define an augments for TTs as mentioned in Appendix D.

This design choice was made to ensure a fair comparison across TTs, as some are not capable of generating diverse transformations from the same text input.

### 2.3.2 One-shot and few-shot based generation

Augmentation resulting from prompt-based generation (PG) refers to the adoption of generative models, specifically LLMs, to generate variations of original MI sessions while preserving contextually relevant information and keeping dialogues coherent with the overall session topic. LLMs represent the prevailing solution for most of language processing tasks, thanks to the versatile management of output diversity through specific parameters, e.g., temperature, top-p sampling, or prompt refinement. Despite these benefits, they can be inconsistent and difficult to steer their generation process toward the required results. Main weaknesses, such as hallucination and poor instruction following, mine the reliability of their replies and might require additional processing layers to filter out irrelevant information.

Unlike TTs, SOTA LLMs support textual inputs with large context windows, enabling our targeted task at the session level by interleaving individual turn talks and interlocutor delimiters. To guarantee high output fidelity and strong instruction following, we adopted several prompting techniques and iteratively refined the initial prompt until satisfied with PG augmentations. Figure 2 in appendix E depicts the final prompting setup, consisting of a system and a user prompt for accurate role assignment and a PG output example with augmented

turn talks. The following formalizes the PG process as TT augmentation one:

$$\phi_{PG}(\mathcal{D}_{train}) = \{(\phi_{PG}^S(S_i), q_i)\}_{i=1}^L \quad (2)$$

where  $\phi_{PG}^S : \mathbb{R}^d \rightarrow \mathbb{R}^d$  denotes a PG sub-function that operates at the session level.

While TTs can guarantee that the resulting output includes an equal number of utterances as the original training set, PG reports higher variance in outcomes. Indeed, due to diverse model parameter sizes, configurations, and underlying functioning, PG may result in shorter or longer sessions compared with the original data.

### 3 Experimental Design

This section outlines the steps for preparing the experimental data and presents the classification architectures, refined promotion strategies, and the experimental pipeline. The targeted objective is addressed by assessing the efficacy of TT and PG augmented data in improving classification performance under the experimental setup below:

- **Data Preparation:** Out of 131 sessions, we selected 97 sessions for the training set, of which 15 are low- and 82 are high-quality MI dialogues. The remaining 34 MI sessions constitute the test set, of which 8 are low- and 26 are high-quality MI sessions. 10 MI sessions (3 low- and 7 high-quality) are removed from the training set to create a validation set further. Such a splitting strategy guarantees that MI quality samples in each set remain representative of the original session distribution.
- **Classification Models:** We employed three classifiers, namely BigBird<sup>3</sup> (Zaheer et al., 2020), Longformer<sup>4</sup> (Beltagy et al., 2020) and ModernBERT<sup>5</sup> (Warner et al., 2024). The rationale behind using BigBird, Longformer, and ModernBERT over usual transformers is their capability to handle long text sequences. Preserving the context of the MI session is crucial to our experimental design, and given that MI sessions are as long as 8000 tokens, these models completely match the experiment’s requirement. The inclusion of ModernBERT

<sup>3</sup>[https://huggingface.co/docs/transformers/en/model\\_doc/big\\_bird](https://huggingface.co/docs/transformers/en/model_doc/big_bird)

<sup>4</sup>[https://huggingface.co/docs/transformers/en/model\\_doc/longformer](https://huggingface.co/docs/transformers/en/model_doc/longformer)

<sup>5</sup>[https://huggingface.co/docs/transformers/en/model\\_doc/ModernBERT](https://huggingface.co/docs/transformers/en/model_doc/ModernBERT)

as the latest modern release of BERT models emphasizes our adoption of SOTA techniques for high-quality outputs. All the experiments adopt accuracy (**Acc**), balanced accuracy (**Bal. Acc.**), F1-score (**F1**), precision (**Prec**), and recall (**Rec**) as evaluation metrics for classification performance of the employed models. The experiment using only the original data  $\mathcal{D}_{train}$ , labeled as AnnMI (Base Data), serves as the baseline to estimate the impact of augmented data in influencing classification performance.

- **Experiment with Domain-trained Model:**

Since our experiments are session and not utterance-based, the classification task becomes complex and challenging, given the complex domain of psychology and a highly skewed and unbalanced distribution of the experimental dataset. Therefore, as a sanity check, we begin by evaluating the performance of the SOTA domain-trained LLM-MentaLLaMA<sup>6</sup> (Yang et al., 2024). This experiment is set to evaluate whether domain-trained models can outperform contemporary classifiers fine-tuned on experimental data, highlighting in favor or against the need for data augmentation.

- **Experiments with LMs:** To evaluate the performance of classification models, TT and PG augmenters, we rigorously conduct three sets of experiments:

i) First, we perform the experiments with the  $\mathcal{D}_{train}$  with the three classifiers.

ii) Then  $\mathcal{D}_{train} + \mathcal{D}_{aug,PG}$  and  $\mathcal{D}_{train} + \mathcal{D}_{aug,TT}$  are used as training data with the three classifiers iteratively, where  $\mathcal{D}_{aug,PG}$  and  $\mathcal{D}_{aug,TT}$  are the augmented dataset generated by PG (i.e., LLMs) and TT augmenters respectively. Here  $\mathcal{D}_{aug}$  can be only one augmented dataset for one experiment.

iii) Now, based on the performance of classification models, the best-performing  $\mathcal{D}_{aug}$  datasets from each  $\mathcal{D}_{aug,PG}$  and  $\mathcal{D}_{aug,TT}$  are selected to make two separate training datasets namely  $\mathcal{D}_{train+aug,PG}$  and  $\mathcal{D}_{train+aug,TT}$  representing the combination of original data with data augmented by each method family.

iv) Finally,  $\mathcal{D}_{train+aug,PG}$  and  $\mathcal{D}_{train+aug,TT}$

<sup>6</sup><https://huggingface.co/klyang/MentaLLaMA-chat-7B>

366 are used as training datasets for each of the  
367 three classifiers iteratively. This experiment  
368 aims to observe the generalization capability  
369 of classifiers over diverse datasets and explore  
370 the combined contributions of effective aug-  
371 mented data.

## 372 4 Results and Discussion

373 This section presents the experimental results of  
374 classification models and empirical analyses for  
375 each data augmentation setup, followed by an eval-  
376 uation of LLMs using a novel tailored evaluation  
377 scheme.

### 378 4.1 Performance of domain-trained LLM 379 MentalLLaMA

380 As mentioned in section 3, we begin our experi-  
381 ment with the SOTA domain-trained LLM Mental-  
382 LLaMA (Yang et al., 2024) to observe if domain-  
383 specific models can outperform domain-agnostic  
384 models on specific task and dataset. Interestingly,  
385 MentalLLaMA performed poorly on the task, as  
386 observed over three iterations. The output expla-  
387 nations produced by MentalLLaMA are unrecog-  
388 nizable, and it mostly predicts all the MI sessions  
389 as "LOW," which is the minority class. This lapse  
390 in the performance of the MentalLLaMA being  
391 SOTA model highlights a clear research gap in  
392 low-resource domains such as mental health, and  
393 it raises concerns about the model’s reliability and  
394 readiness for domain-specific NLP downstream  
395 tasks. It also indicates the need to enrich the low-  
396 resource domain to investigate optimal practices  
397 for generating high-quality data, enabling LMs and  
398 LLMs to better interpret the nuances and intricac-  
399 ies of the complex domain, leading to reliable  
400 performance.

### 401 4.2 Classifiers performance with PG-based 402 augmented datasets

403 The performance of the three classifiers with each  
404 of the 13 LLM-augmented datasets is summarized  
405 in Table 1. For ease of interpretation, the best-  
406 performing augmented dataset and models are high-  
407 lighted in bold. Table 1 shows that most of the  
408 LLM-augmented datasets have improved the em-  
409 ployed LMs’ accuracy compared to the baseline.  
410 For instance, BigBird has improved by **5.88%** with  
411 Llama 3.2, Longformer by **17.65%** with Smollm2,  
412 and ModernBERT by **2.94%** with Nemotron-Mini.  
413 It is also noteworthy that the baseline performance  
414 of ModernBERT is significantly higher and the

415 highest among the three LMs, which emphasizes  
416 the efficient architecture of ModernBERT, designed  
417 to reduce bias and pay particular attention to effi-  
418 ciency. Also, peculiarly, Longformer has the lowest  
419 baseline performance, but augmentation shines and  
420 improves its performance by a noteworthy **17.65%**.

### 421 4.3 Classifiers Performance with TT-based 422 augmented datasets

423 The performance of each of the three classifiers  
424 with each of the 5 NLP-augmented datasets is  
425 summed up in Table 3 in Appendix G. The best-  
426 performing data and models are highlighted in  
427 bold. As evident from Table 3, most of the NLP-  
428 augmented datasets have also improved the LMs’  
429 accuracy compared to the baseline. For instance,  
430 BigBird has improved by **5.88%** with Backtransla-  
431 tion, Longformer by **17.65%** with ProtAugment  
432 and Tweet style, and ModernBERT by **2.24%** with  
433 Backtranslation. Besides the best performance, the  
434 "Basic style" is consistently the lowest performer.  
435 Additionally, in retrospect, Longformer has the  
436 lowest baseline performance, and augmentation  
437 has been very useful, yielding an improvement of  
438 **17.65%**.

### 439 4.4 Performance of LMs with mixed 440 augmented dataset

441 To evaluate the generalization capability of LMs  
442 over diversified data, i.e., the training dataset  
443 containing mixed-style synthetic data, we used  
444 a combination for LLM- and NLP-augmented  
445 datasets. From LLM, we have Nemotron-Mini(4B)  
446 and Phi4(14B), the best performers averaged over  
447 all experiments. Similarly, from NLP-augmented  
448 datasets, we select Backtranslation, ProtAugment  
449 and Tweet style for the best performance (averaged  
450 over all experiments). The performance of each  
451 of the three classifiers individually with the mixed  
452 LLM- and NLP-augmented dataset is presented in  
453 Table 2 and the best-performing data and models  
454 are highlighted in bold. The LMs BigBird and  
455 Longformer have improved their performance with  
456 the mixed, augmented dataset. To be specific,  
457 BigBird has observed an increase of **3.64%** in  
458 accuracy with combined LLM-augmented data  
459 and Longformer has scored a **17.65%** increase  
460 in accuracy with LLM-augmented data. As  
461 mentioned before, ModernBERT, being SOTA  
462 and a robust model, has set the highest baseline  
463 score among the LMs. Therefore, it is imperative  
464 to fine-tune ModernBERT on high-quality data to

Model	Augmented Data	Acc	Bal. Acc.	F1	Prec	Rec
<b>BigBird</b>	AnnoMI (Base Data)	0.8824	<b>0.9231</b>	0.8892	0.9216	0.8824
	ChatGPT (4.0)	0.7647	0.5000	0.6627	0.5848	0.7647
	DeepSeek (v3)	0.7059	0.5048	0.6674	0.6451	0.7059
	Gemini-Flash	0.8824	0.7500	0.8669	0.8980	0.8824
	Llama3.1 (70B)	0.7941	0.6923	0.7892	0.7859	0.7941
	Llama3.2 (3B)	<b>0.9412</b>	0.9183	<b>0.9412</b>	<b>0.9412</b>	<b>0.9412</b>
	Llama3.3 (70B)	0.8529	0.6875	0.8260	0.8767	0.8529
	Mistral-Nemo (12B)	0.8235	0.6250	0.7797	0.8566	0.8235
	Nemotron-Mini (4B)	0.8824	0.7933	0.8761	0.8789	0.8824
	Phi4 (14B)	0.9118	0.8125	0.9040	0.9209	0.9118
	Qwen2.5 (32B)	0.8235	0.6683	0.8004	0.8137	0.8235
	Qwen2.5 (72B)	0.7647	0.5000	0.6627	0.5848	0.7647
	Smollm2 (1.7B)	0.7941	0.5625	0.7263	0.8378	0.7941
<b>Longformer</b>	AnnoMI (Base Data)	0.7647	0.5000	0.6627	0.5848	0.7647
	ChatGPT (4.0)	0.7647	0.5000	0.6627	0.5848	0.7647
	DeepSeek (v3)	0.8529	0.7308	0.8400	0.8475	0.8529
	Gemini-Flash	0.7647	0.5000	0.6627	0.5848	0.7647
	Llama3.1 (70B)	0.7647	0.5000	0.6627	0.5848	0.7647
	Llama3.2 (3B)	0.7647	0.5000	0.6627	0.5848	0.7647
	Llama3.3 (70B)	0.7647	0.5000	0.6627	0.5848	0.7647
	Mistral-Nemo (12B)	0.8824	0.7500	0.8669	0.8980	0.8824
	Nemotron-Mini (4B)	0.7941	0.8221	0.8077	0.8550	0.7941
	Phi4 (14B)	0.8529	0.7308	0.8400	0.8475	0.8529
	Qwen2.5 (32B)	0.8235	0.6683	0.8004	0.8137	0.8235
	Qwen2.5 (72B)	0.8235	0.7548	0.8235	0.8235	0.8235
	Smollm2 (1.7B)	<b>0.9412</b>	<b>0.8750</b>	<b>0.9381</b>	<b>0.9454</b>	<b>0.9412</b>
<b>ModernBERT</b>	AnnoMI (Base Data)	0.9118	0.8125	0.9040	0.9209	0.9118
	ChatGPT (4.0)	0.8529	0.6875	0.8260	0.8767	0.8529
	DeepSeek (v3)	0.8824	0.8365	0.8824	0.8824	0.8824
	Gemini-Flash	0.7941	0.6490	0.7760	0.7740	0.7941
	Llama3.1 (70B)	0.8529	0.7740	0.8494	0.8478	0.8529
	Llama3.2 (3B)	0.7647	0.5865	0.7339	0.7294	0.7647
	Llama3.3 (70B)	0.7941	0.6058	0.7564	0.7736	0.7941
	Mistral-Nemo (12B)	0.8529	0.8173	0.8558	0.8604	0.8529
	Nemotron-Mini (4B)	<b>0.9412</b>	<b>0.8750</b>	<b>0.9381</b>	<b>0.9454</b>	<b>0.9412</b>
	Phi4 (14B)	0.9118	0.8125	0.9040	0.9209	0.9118
	Qwen2.5 (32B)	0.8235	0.6683	0.8004	0.8137	0.8235
	Qwen2.5 (72B)	0.8824	0.7500	0.8669	0.8980	0.8824
	Smollm2 (1.7B)	0.8529	0.7308	0.8400	0.8475	0.8529

Table 1: Classification Metrics for different models and LLM-generated augmented data.

465	further improve performance on mixed training	512
466	data.	513
467		514
468	Across experiments, session-level PG consis-	515
469	tently improves balanced accuracy over utterance-	
470	level TT, particularly for minority-class MI ses-	
471	sions, indicating that long-range discourse is criti-	
472	cal for MI quality modeling.	
473	<b>5 Evaluation Scheme and</b>	
474	<b>LLM-generated Data Analysis</b>	
475	This section presents the evaluation scheme we	
476	developed, centered on the experimental setup and	
477	dataset, and how we implemented it, followed by a	
478	detailed explanation of each LLM’s performance.	
479	<b>5.1 Evaluation scheme development</b>	
480	For the quantitative and qualitative evaluation of	
481	LLMs, we developed a novel evaluation scheme	
482	with the help of a domain expert. In the evaluation	
483	scheme, we parameterize the attributes of MI ses-	
484	sions to capture the full spectrum of their quality	
485	and content. The evaluation scheme is applied to	
486	the representative sample taken from $\mathcal{D}_{train}$ and	
487	each of the $\mathcal{D}_{augllm}$ containing original $MI_{orig.}$	
488	and synthetic MI sessions $MI_{syn.}$ . The outcome is	
489	the average score over the representative sample.	
490		
491	• <b>Semantic Drift:</b> It refers to the change in	
492	the meaning of words or phrases over time or	
493	across contexts. In our case, we use it to indi-	
494	cate how much semantic deviation happened	
495	in $MI_{syn.}$ with respect to (w.r.t) to $MI_{orig.}$ .	
496	It is measured on a scale of High, Moderate,	
497	and Low.	
498		
499	• <b>Hallucination:</b> It refers to the false or un-	
500	supported information generation w.r.t to the	
501	input data. We use it to observe how often	
502	$MI_{syn.}$ has made the wrong interpretation of	
503	the therapeutic notion w.r.t $MI_{orig.}$ . It is mea-	
504	sured on a scale of High, Moderate, and Low.	
505		
506	• <b>MI Session Length:</b> We use this parameter	
507	to observe if extreme changes in the length	
508	of $MI_{syn.}$ w.r.t $MI_{orig.}$ are significantly con-	
509	tributing to the quality and false semantic de-	
510	viation in $MI_{syn.}$ . It is measured in Longer,	
511	Equal to Shorter w.r.t to $MI_{orig.}$ .	
	• <b>Prompt Compliance:</b> We use this parameter	
	to observe how strictly an LLM adheres to	
	the instructions in the prompts. It signifies	
	user-friendliness, input understanding, and de-	
	sirable output format on the LLM level. It is	
	measured on a scale of High, Moderate, and	
	Low.	
	• <b>Fluency and Emotional Understanding:</b> We	
	use it to refer to the naturalness, coherency,	
	emotional stance, and smoothness of lan-	
	guage in $MI_{syn.}$ w.r.t $MI_{orig.}$ . The emotional	
	stance here indicates if the emotion the client	
	or therapist displays in $MI_{orig.}$ is presented	
	more strongly or mildly in $MI_{syn.}$ . For in-	
	stance, if the client expresses mild tiredness	
	in $MI_{orig.}$ while in $MI_{syn.}$ if perceived as ex-	
	tremely tired and sick, indicates <b>Low</b> fluency	
	and emotional understanding.	
	<b>5.2 LLM-generated data analysis</b>	
	The performance of each LLM and the LLM-	
	augmented data over the five parameters is pre-	
	sented in Table 4 placed in Appendix H, and we	
	have consolidated the point-to-point discussion as	
	follows.	
	• <b>Semantic Drift:</b> Interestingly, the LLMs	
	have skipped extreme semantic drift (except	
	Smollm2) and have only shown Moderate to	
	Low semantic deviation.	
	• <b>Hallucination:</b> None of the LLMs have	
	shown extreme hallucination. However, we	
	noticed different hallucination variations in	
	the augmented MI sessions. For instance,	
	Gemini introduced a pseudo-name for the	
	client to make the conversation flow more nat-	
	urally. However, it does not change the con-	
	text; what it does change is the client’s gen-	
	der, which might conflict with specific therapy	
	courses.	
	• <b>MI Session Length:</b> In this case, LLMs	
	have shown all variations on the scale, i.e.,	
	while Nemotron and Smollm2 have produced	
	mostly extra-long MI sessions, the remaining	
	LLMs have generated shorter sessions. Fur-	
	ther, to give a fine-grain quantitative overview,	
	we have presented the distribution of utter-	
	ances per $MI_{orig.}$ vs. $MI_{syn.}$ for each LLM	
	separately, as mentioned in Figure 3, in Ap-	
	pendix F.	
	• <b>Prompt Compliance:</b> None of the LLMs	
	showed complete ignorance or inability to	
	follow the prompt instruction. While some	

Model	Augmented Data	Acc	Bal. Acc.	F1	Prec	Rec
<b>BigBird</b>	AnnoMI (Base Data)	0.8824	<b>0.9231</b>	0.8892	<b>0.9216</b>	0.8824
	Backtranslation + ProtAugment + Tweet Style	0.8824	0.8798	0.8865	0.8975	0.8824
	Nemotron-Mini (4B) + Phi-4 (14B)	<b>0.9118</b>	0.8558	<b>0.9097</b>	0.9097	<b>0.9118</b>
<b>Longformer</b>	AnnoMI (Base Data)	0.7647	0.5000	0.6627	0.5848	0.7647
	Backtranslation + ProtAugment + Tweet Style	<b>0.9412</b>	0.8750	<b>0.9381</b>	<b>0.9454</b>	<b>0.9412</b>
	Nemotron-Mini (4B) + Phi-4 (14B)	0.8824	<b>0.8798</b>	0.8865	0.8975	0.8824
<b>ModernBERT</b>	AnnoMI (Base Data)	<b>0.9118</b>	<b>0.8125</b>	<b>0.9040</b>	<b>0.9209</b>	<b>0.9118</b>
	Backtranslation + ProtAugment + Tweet Style	0.8235	0.6683	0.8004	0.8137	0.8235
	Nemotron-Mini (4B) + Phi-4 (14B)	0.8235	0.6683	0.8004	0.8137	0.8235

Table 2: Performance of each LM with combined LLM- and NLP-augmented datasets.

LLMs are seamless, a few exhibit moderate difficulties in interpreting or following the chain of instructions. The most common issue is the inability to produce the desired outcome in the prescribed format for the same prompts. And Gemini, Llama, Mistral Nemo, Phi4, Qwen, and Smollm2 were prone to it.

- **Fluency and Emotional Understanding:** While all the LLMs have shown great fluency, capturing the actual depth of emotions and context is somewhat more challenging. For instance, DeepSeek, Gemini and LLama have been observed to elicit stronger emotional responses. Table 5 in Appendix I provides a deeper look into each LLM’s performance and reflects the extent of MI adherence in  $MI_{syn.}$ , summarizing the main trends and observations about each LLM.

## 6 Conclusion and Future Work

In this study, we investigate the performance of 13 LLM for in-context, human-like MI dialogue generation at the entire conversation level, of which **DeepSeek** is the latest release and SOTA. We also investigate TT-based augmentation techniques applied at the utterance level to assess whether the context and intricacies of complex domains can still be captured at the utterance level. Further, to analyze the quality of LLM and TT-generated augmented data, we used BigBird, Longformer, and ModernBERT classifiers designed to process long sequences, of which ModernBERT is the latest release and SOTA transformer model. Our results show for such a difficult task to perform classification over considering an entire conversation on a severely unbalanced and small dataset belonging

to a complex domain, the classifiers have demonstrated exceptional performance, with the highest balanced score being **98.08%**. Then, we employ our meticulously developed evaluation scheme parametrized to measure semantic drift, hallucination, MI adherence, prompt compliance, fluency, context, and emotional interpretation in LLMs and facilitate fine-grained qualitative and quantitative analysis. As a resource contribution, we provide the first-ever high-quality dataset containing 1764 low- and high-quality MI dialogues and 81k talk turns between therapist and client. Based on our empirical analysis, ChatGPT, DeepSeek, and Gemini have shown impressive capabilities in using open-ended questions to encourage exploration in MI, reflecting feelings, and affirming emotions to build rapport, encouraging small, manageable steps instead of overwhelming directives, and fostering collaboration and autonomy, allowing the client to choose changes and that may lead to change talk. For future work, we aim to explore additional domain-specific LLMs, investigate possibilities for fine-tuning LLMs on our data for domain-specific NLP downstream tasks (automatic MI session annotation, summary generation, etc.), and further refine our evaluation framework to extend its applicability to related areas of mental health research.

## Limitations

While this work presents a comprehensive benchmark for session-level LLM-based data augmentation in Motivational Interviewing, several limitations should be acknowledged. First, although the curated dataset is the largest session-level MI resource used for this type of study to date, it remains modest in scale compared to general-domain dialogue corpora. The dataset also exhibits class im-

631 balance that reflects real-world MI practice, which  
632 may influence classifier performance and limit gen-  
633 eralization to settings with different label distribu-  
634 tions.

635 Second, the MI sessions used in this study are  
636 drawn from a specific clinical and cultural context.  
637 As MI practices can vary across populations, lan-  
638 guages, and institutional settings, the findings may  
639 not directly generalize to other forms of counseling  
640 or to non-English MI interactions without further  
641 validation.

642 Third, our data augmentation strategies rely on  
643 contemporary large language models whose out-  
644 puts are inherently sensitive to prompt design, de-  
645 coding parameters, and model updates. While  
646 we employ standardized prompts and controlled  
647 generation settings, different prompt formulations  
648 or future model versions may yield varying re-  
649 sults. In addition, despite incorporating qualitative  
650 checks for hallucination and semantic drift, LLM-  
651 generated data may still introduce subtle artifacts  
652 that are not fully captured by automated evaluation  
653 metrics.

654 Fourth, our evaluation focuses on downstream  
655 MI session classification rather than clinical out-  
656 comes or human-centered measures of counseling  
657 effectiveness. Consequently, improvements in clas-  
658 sification performance should not be interpreted as  
659 direct evidence of improved therapeutic quality or  
660 clinical utility. Any real-world deployment of such  
661 systems would require rigorous human-in-the-loop  
662 validation and ethical oversight.

663 Finally, while we examine multiple long-context  
664 classification models, our study does not exhaus-  
665 tively explore architectural variations or fine-tuning  
666 strategies. Future work could investigate alternative  
667 model designs, multilingual settings, and tighter in-  
668 tegration of domain expert feedback during both  
669 data generation and evaluation.

## 670 References

671 Nicklaus Badyal, Derek Jacoby, and Yvonne Coady.  
672 2023. [Intentional biases in LLM responses](#). In *14th*  
673 *IEEE Annual Ubiquitous Computing, Electronics &*  
674 *Mobile Communication Conference, UEMCON 2023,*  
675 *New York, NY, USA, October 12-14, 2023*, pages 502–  
676 506. IEEE.

677 Simone Balloccu, Ehud Reiter, Karen Jia-Hui Li, Rafael  
678 Sargsyan, Vivek Kumar, Diego Reforgiato, Daniele  
679 Riboni, and Ondrej Dusek. 2024a. [Ask the experts:](#)  
680 [sourcing a high-quality nutrition counseling dataset](#)  
681 [through human-AI collaboration](#). In *Findings of the*

*Association for Computational Linguistics: EMNLP*  
2024, pages 11519–11545, Miami, Florida, USA.  
Association for Computational Linguistics.

685 Simone Balloccu, Patrícia Schmidtová, Mateusz Lango,  
686 and Ondrej Dusek. 2024b. [Leak, cheat, repeat: Data](#)  
687 [contamination and evaluation malpractices in closed-](#)  
688 [source LLMs](#). In *Proceedings of the 18th Confer-*  
689 *ence of the European Chapter of the Association*  
690 *for Computational Linguistics (Volume 1: Long Pa-*  
691 *pers)*, pages 67–93, St. Julian’s, Malta. Association  
692 for Computational Linguistics.

693 Erkan Basar, Iris Hendrickx, Emiel Kraemer, Gert-Jan  
694 Bruijn, and Tibor Bosse. 2024. [To what extent are](#)  
695 [large language models capable of generating substan-](#)  
696 [tial reflections for motivational interviewing counsel-](#)  
697 [ing chatbots? a human evaluation](#). In *Proceedings of*  
698 *the 1st Human-Centered Large Language Modeling*  
699 *Workshop*, pages 41–52, TBD. ACL.

700 Erkan Basar, Xin Sun, Iris Hendrickx, Jan de Wit, Tibor  
701 Bosse, Gert-Jan De Bruijn, Jos A. Bosch, and Emiel  
702 Kraemer. 2025. [How well can large language models](#)  
703 [reflect? a human evaluation of LLM-generated re-](#)  
704 [flections for motivational interviewing dialogues](#). In  
705 *Proceedings of the 31st International Conference on*  
706 *Computational Linguistics*, pages 1964–1982, Abu  
707 Dhabi, UAE. Association for Computational Linguistics.  
708

709 Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020.  
710 [Longformer: The long-document transformer](#).

711 Emily M. Bender, Timnit Gebru, Angelina McMillan-  
712 Major, and Shmargaret Shmitchell. 2021. [On the](#)  
713 [dangers of stochastic parrots: Can language mod-](#)  
714 [els be too big?](#) In *Proceedings of the 2021 ACM*  
715 *Conference on Fairness, Accountability, and Trans-*  
716 *parency, FAccT ’21*, page 610–623, New York, NY,  
717 USA. Association for Computing Machinery.

718 Andrew Brown, Ash Tanuj Kumar, Osnat Melamed,  
719 Imtihan Ahmed, Yu Hao Wang, Arnaud Deza, Marc  
720 Morcos, Leon Zhu, Marta Maslej, Nadia Minian,  
721 Vidya Sujaya, Jodi Wolff, Olivia Doggett, Mathew  
722 Iantorno, Matt Ratto, Peter Selby, and Jonathan Rose.  
723 2023. [A motivational interviewing chatbot with gen-](#)  
724 [erative reflections for increasing readiness to quit](#)  
725 [smoking: Iterative development study](#). *JMIR Ment*  
726 *Health*, 10:e49132.

727 Andrew Brown, Jiading Zhu, Mohamed Abdelwahab,  
728 Alec Dong, Cindy Wang, and Jonathan Rose. 2024.  
729 [Generation, distillation and evaluation of motiva-](#)  
730 [tional interviewing-style reflections with a founda-](#)  
731 [tional language model](#). In *Proceedings of the 18th*  
732 *Conference of the European Chapter of the Associa-*  
733 *tion for Computational Linguistics (Volume 1: Long*  
734 *Papers)*, pages 1241–1252, St. Julian’s, Malta. Asso-  
735 ciation for Computational Linguistics.

736 Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
737 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
738 Neelakantan, Pranav Shyam, Girish Sastry, Amanda

739	Askeel, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. <a href="#">Language models are few-shot learners</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 1877–1901. Curran Associates, Inc.	796
740		797
741		798
742		799
743		800
744		801
745		802
746	Mingda Chen, Jingfei Du, Ramakanth Pasunuru, Todor Mihaylov, Sridi Iyer, Veselin Stoyanov, and Zornitsa Kozareva. 2022. <a href="#">Improving in-context few-shot learning via self-supervised training</a> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3558–3573, Seattle, United States. Association for Computational Linguistics.	803
747		804
748		805
749		806
750		807
751		808
752		809
753		810
754		811
755		812
756	Dorottya Demszky, Diyi Yang, David S Yeager, Christopher J Bryan, Margaret Clapper, Susannah Chandhok, Johannes C Eichstaedt, Cameron Hecht, Jeremy Jamieson, Meghann Johnson, and 1 others. 2023. Using large language models in psychology. <i>Nature Reviews Psychology</i> , 2(11):688–701.	813
757		814
758		815
759		816
760		817
761		818
762		819
763		
764		
765		
766		
767		
768		
769		
770		
771		
772		
773		
774		
775		
776		
777		
778		
779		
780		
781		
782		
783		
784		
785		
786		
787		
788		
789		
790		
791		
792		
793		
794		
795		
796		
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		
810		
811		
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	<i>the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 5609–5621, Torino, Italia. ELRA and ICCL.	908
852		909
853		910
854		911
855	Ruixiang Tang, Xiaotian Han, Xiaoqian Jiang, and Xia Hu. 2023. <i>Does synthetic data generation of llms help clinical text mining?</i> <i>CoRR</i> , abs/2303.04360.	912
856		913
857		914
858	Yuqing Wang, Yun Zhao, and Linda Petzold. 2023. Are large language models ready for healthcare? a comparative study on clinical language understanding. In <i>Proceedings of the 8th Machine Learning for Healthcare Conference</i> , volume 219 of <i>Proceedings of Machine Learning Research</i> , pages 804–823. PMLR.	915
859		916
860		917
861		918
862		919
863		
864	Benjamin Warner, Antoine Chaffin, Benjamin Clavié, Orion Weller, Oskar Hallström, Said Taghadouini, Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom Aarsen, Nathan Cooper, Griffin Adams, Jeremy Howard, and Iacopo Poli. 2024. <i>Smarter, better, faster, longer: A modern bidirectional encoder for fast, memory efficient, and long context finetuning and inference</i> .	920
865		921
866		922
867		923
868		924
869		925
870		926
871		927
872	Anuradha Welivita and Pearl Pu. 2022. <i>Curating a large-scale motivational interviewing dataset using peer support forums</i> . In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 3315–3330, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.	928
873		929
874		930
875		931
876		932
877		933
878	Yuqi Wu, Kaining Mao, Yanbo Zhang, and Jie Chen. 2024. <i>Callm: Enhancing clinical interview analysis through data augmentation with large language models</i> . <i>IEEE Journal of Biomedical and Health Informatics</i> , 28(12):7531–7542.	934
879		
880		
881		
882		
883	Zixiu Wu, Simone Balloccu, Vivek Kumar, Rim Helaoui, Diego Reforgiato Recupero, and Daniele Riboni. 2023. <i>Creation, analysis and evaluation of annomi, a dataset of expert-annotated counselling dialogues</i> . <i>Future Internet</i> , 15(3).	935
884		936
885		937
886		938
887		939
888	Zixiu Wu, Simone Balloccu, Vivek Kumar, Rim Helaoui, Ehud Reiter, Diego Reforgiato Recupero, and Daniele Riboni. 2022. <i>Anno-mi: A dataset of expert-annotated counselling dialogues</i> . In <i>ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 6177–6181.	940
889		941
890		942
891		943
892		944
893		945
894		946
895	Xuhai Xu, Bingsheng Yao, Yuanzhe Dong, Saadia Gabriel, Hong Yu, James Hendler, Marzyeh Ghassemi, Anind K. Dey, and Dakuo Wang. 2024. <i>Mental-llm: Leveraging large language models for mental health prediction via online text data</i> . <i>Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.</i> , 8(1).	947
896		948
897		949
898		950
899		951
900		952
901	Kailai Yang, Tianlin Zhang, Ziyang Kuang, Qianqian Xie, Jimin Huang, and Sophia Ananiadou. 2024. <i>Mental-lama: Interpretable mental health analysis on social media with large language models</i> . In <i>Proceedings of the ACM Web Conference 2024, WWW '24</i> , page 4489–4500, New York, NY, USA. Association for Computing Machinery.	953
902		954
903		955
904		956
905		957
906		958
907		
	Jiayi Yuan, Ruixiang Tang, Xiaoqian Jiang, and Xia Hu. 2023. <i>Large language models for healthcare data augmentation: An example on patient-trial matching</i> . <i>AMIA Annu. Symp. Proc.</i> , 2023:1324–1333.	
	Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. 2020. <i>Big bird: Transformers for longer sequences</i> . In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 17283–17297. Curran Associates, Inc.	
	<b>Ethical Considerations</b>	
	This work focuses on methodological evaluation and benchmarking rather than deployment in real-world clinical settings. All MI data used in this study are de-identified and handled in accordance with applicable data protection guidelines. While our use of large language models enables scalable data augmentation, such models may generate hallucinated or misleading content, particularly in sensitive mental health contexts. Accordingly, the generated data and trained models should not be used for diagnostic or therapeutic purposes without rigorous human oversight. Any future application of these methods in practice would require careful validation, domain expert involvement, and ethical review to ensure patient safety and responsible use.	
	<b>A Related Work</b>	
	The research field of integrating LLMs into therapeutic settings, particularly MI, is still evolving, and there is not much work done to answer specific research questions (Wu et al., 2024; Yuan et al., 2023; Li et al., 2023b; Tang et al., 2023), such as data augmentation at the MI session level, using LLMs for mass annotation tasks, etc. In fact, since 2016, only a few works (Pérez-Rosas et al., 2016; Welivita and Pu, 2022; Kumar et al., 2023; Wu et al., 2023) have contributed to the development of public datasets, as mentioned in (Sun et al., 2024). (Balloccu et al., 2024a) can be considered a holistic work to generate the HAI-Coaching dataset by collecting approximately 2,400 real-world dietary struggles and around 97,000 supportive responses using LLMs, evaluated by nutrition experts for safety and quality. Among a few available works, (Basar et al., 2025) evaluates the ability of LLMs to generate reflective listening responses in MI dialogues. The study assesses whether LLM-generated responses align with human therapists' reflections and effectiveness in therapeutic conversations. Early credits can be given to (Brown	

959 [et al., 2023](#)) to invoke MI-based conversation to  
960 help smokers move toward quitting smoking. In  
961 the subsequent work, ([Brown et al., 2024](#)) inves-  
962 tigate how LLMs can be fine-tuned to generate  
963 MI-style reflections for use in therapeutic chatbots.  
964 It evaluates the effectiveness of AI-generated reflec-  
965 tions in maintaining engagement and supporting  
966 behavior change. The work in ([Basar et al., 2024](#))  
967 focuses on the quality, coherence, and therapeu-  
968 tic value of AI-generated reflections compared to  
969 human-generated ones. It provides insights into  
970 how LLMs can generate meaningful reflections in  
971 counseling. As evident from the above works, there  
972 are several gaps in the research of LLM integration  
973 in effective and reliable healthcare services. There-  
974 fore, with this work, we aim to provide a clear path  
975 to integrate LLMs in the therapeutic scenario for  
976 benchmarking, data augmentation, automatic eval-  
977 uation of MI sessions, using SOTA, and further  
978 improvement.

979 **B High- and Low-quality MI Sessions**

980 **C Risk Minimization Problem**

981 **D Defining An Augmenter for TTs**

982 **E Prompt-based Generation Example**

983 **F Distribution of Talk-turns in**  
984 **LLM-LLM-augmented datasets**

985 **G Classification Metrics For Different**  
986 **Models and LLM-generated**  
987 **Augmented Data.**

988 **H Evaluation Scheme MI Dialogues**  
989 **Generated by LLMs**

990 **I Fine-grained Evaluation of LLMs**  
991 **Performance**

MI Dialogue (High-quality)		MI Dialogue (Low-quality)	
<b>Therapist</b>	What else has brought you here?	<b>Therapist</b>	Hi. You are?
<b>Client</b>	They found a little bag of weed.	<b>Client</b>	Yeah.
<b>Therapist</b>	Mm-hmm.	<b>Therapist</b>	Okay. Um, what brings you in?
<b>Client</b>	It's not like I was gonna sell it or anything. It's like	<b>Client</b>	Well, I'm supposed to come to learn about my--
<b>Therapist</b>	Just enough for your own use?	<b>Therapist</b>	how to manage my diabetes. Just got diagnosed.
<b>Client</b>	Yeah, whatever, you know. It's California, no one	<b>Therapist</b>	Mm-hmm. I see- I see that your A1C or your
	cares.		blood glucose is very high.
<b>Therapist</b>	Yeah. Yeah. Yeah.	... ..	... ..
<b>Client</b>	How long is this gonna take?		That's, uh, that's a huge amount of stuff to do and,
... ..	... ..	<b>Client</b>	I'm not- I'm not certain I can do all that, and on
	Okay, mm-hmm. Well, I'm wondering then if it's		top of it, it's-it's-- if I don't get it all done, all this
	okay with you and-and I have a form that we can use	<b>Therapist</b>	stuff happens to me?
	if you'd like, um, just to put out some different ideas		So, it sure can.
	on an agenda that we could share as we're working	<b>Client</b>	[sighs] So, uh, you know, we-we need to start
<b>Therapist</b>	together so that I-I can support you to be successful		looking at doing all those different things to get
	in getting this behind you. Would-would this feel	<b>Therapist</b>	your diabetes under control 'cause your-your
	like a natural time to-to move towards figuring out		number is very high.
	where we might go together?	<b>Client</b>	Well, I-I

Figure 1: A sample of a). high-quality and b) low-quality MI sessions.

$$\begin{aligned}
\theta^*(\phi) &= \arg \min_{\theta} \sum_{(S,q) \in \mathcal{D}_{train} \cup \phi(\mathcal{D}_{train})} \mathcal{L}(f_{\theta}(\Gamma(S)), q) \\
\phi^* &= \arg \min_{\phi \in \Phi} \sum_{(S,q) \in \mathcal{D}_{test}} |f_{\theta^*(\phi)}(\Gamma(S)) - q|
\end{aligned} \tag{3}$$

Where  $\mathcal{L}(\cdot, \cdot)$  is a training loss function for classification tasks, e.g., binary cross-entropy. In other words, our objective is to identify the augmenter that can raise the bar in MI quality classification by minimizing the misclassification error on a fixed testing set.

$$\begin{aligned}
\phi_{TT}^S(S_i) &= \phi_{TT}^S(t_{i1}, t_{i2}, \dots, t_{iM_i}) = (\phi_{TT}^t(t_{i1}), \phi_{TT}^t(t_{i2}), \dots, \phi_{TT}^t(t_{iM_i})), \\
\phi_{TT}(\mathcal{D}_{train}) &= \{(\phi_{TT}^S(S_i), q_i)\}_{i=1}^L = \{(\phi_{TT}^t(t_{i1}), \dots, \phi_{TT}^t(t_{iM_i})), q_i\}_{i=1}^L
\end{aligned} \tag{4}$$

Where  $\phi_{TT}^t : \mathbb{R}^e \rightarrow \mathbb{R}^e$  denotes a TT sub-function that operates at the talk level. Note that the TTs adopted at the utterance level imply  $|\mathcal{D}_{Aug}| = |\mathcal{D}_{train}|$ , unless  $\phi_{TT}^t$  is repeatedly applied to the same talk. In other words, we constrained TTs augmenters to generate the same amount of utterances included in the original training set, maintaining also the same order.

**System Prompt**

You are an expert psychologist specializing in motivational interviewing with exceptional rephrasing abilities and a deep understanding of therapeutic techniques.  
Your task is to rephrase each line of a therapy session to enhance clarity, empathy, and professional tone without altering the core meaning or intent of the conversation.

**Task Instructions:**

- Review the context of each interaction between the client and therapist.
- Identify implicit emotional or cognitive themes in the client's responses.
- Rephrase each line to improve tone, empathy, or clarity, reflecting motivational interviewing principles, including open-ended questions, reflective listening, and affirmation.
- Do not create any new content or expand on any ideas.

**Chain of Thought Process:**

- Step 1: Carefully review each client and therapist utterance.
- Step 2: Interpret the client's underlying emotions or cognitive themes.
- Step 3: Rephrase each response thoughtfully, ensuring it aligns with the original meaning while subtly improving the dialogue.
- Step 4: Preserve the original conversational flow and structure without adding new lines or content.

**Output Format:**

- Respond only with rephrased dialogue lines in the format:  
Therapist: [rephrased therapist response]  
Client: [rephrased client response]
- Avoid any narrative descriptions, explanations, or additional text outside of the direct dialogue lines.

---

**User Prompt**

Please rephrase the following therapy session, focusing on improving clarity, empathy, and tone while preserving the original meaning and flow.

Here is the session to rephrase:

*Client: I've been feeling like my wife has been on my case a lot lately about improving my health habits, but to be honest, I don't really see what the issue is.*  
*Therapist: Just to make sure I understand, you're feeling a bit puzzled by your wife's concerns about your health habits?*  
*Client: Yeah, I feel fine. I mean, my job is physically demanding, so I'm always on the go and sweating out anything unhealthy I might eat.*  
*Therapist: That makes sense. It sounds like you're relying on your physical activity to balance out your diet.*

Provide only rephrased dialogue in the format: Sample MI Session  
 Therapist: [rephrased therapist response]  
 Client: [rephrased client response]

Do not add new dialogue or text outside the provided lines.

**LLM Output**

**Therapist:** Just to clarify, it sounds like you're feeling uncertain about why your wife is so concerned about your health habits  
**Client:** Yeah, I feel fine. My job keeps me physically active, so I figure I'm burning off anything unhealthy I eat.  
**Therapist:** I see. It sounds like you rely on your level of physical activity to counterbalance your diet.

Figure 2: A structured chat interface showcasing a system prompt, user prompt, and LLM output in a MI rephrasing task. The system prompt instructs the model to enhance clarity, empathy, and professional tone while preserving meaning. The user provides a sample therapy session, and the LLM generates a refined version, demonstrating its ability to improve conversational quality while maintaining fidelity to the original dialogue.

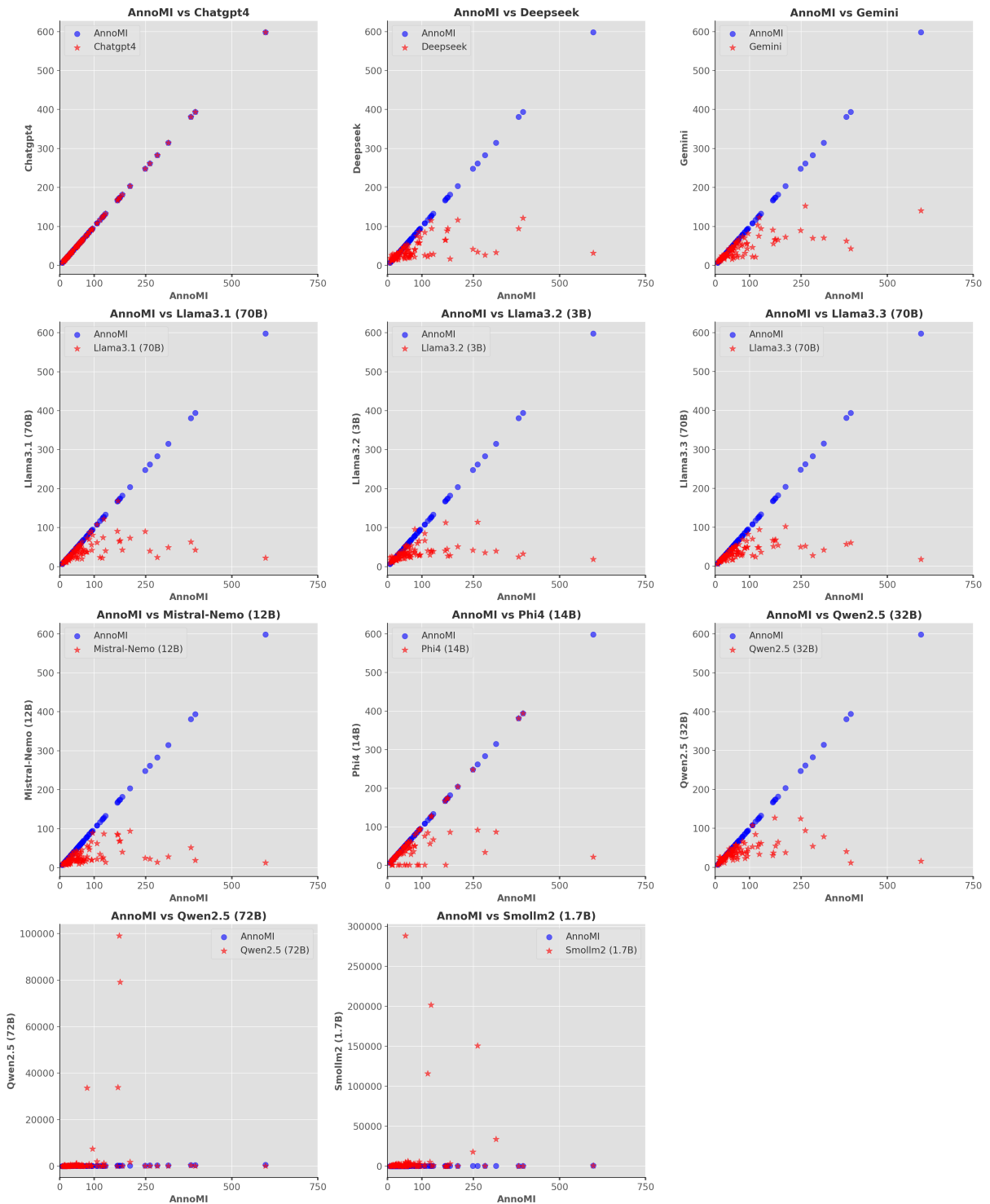


Figure 3: Distribution of the number of turn talks for each  $MI_{orig.}$  vs  $MI_{syn.}$  for all LLM-augmented datasets.

Model	Augmented Data	Acc	Bal. Acc.	F1	Prec	Rec
<b>BigBird</b>	AnnoMI (Base Data)	0.8824	0.9231	0.8892	0.9216	0.8824
	Backtranslation	<b>0.9706</b>	<b>0.9808</b>	<b>0.9712</b>	<b>0.9739</b>	<b>0.9706</b>
	Basic style	0.8529	0.7308	0.8400	0.8475	0.8529
	ProtAugment	0.9412	0.8750	0.9381	0.9454	0.9412
	Synonym substitution	0.8235	0.6250	0.7797	0.8566	0.8235
	Tweet style	0.9118	0.8558	0.9097	0.9097	0.9118
<b>Longformer</b>	AnnoMI (Base Data)	0.7647	0.5000	0.6627	0.5848	0.7647
	Backtranslation	0.8824	0.8798	0.8865	0.8975	0.8824
	Basic style	0.7941	0.5625	0.7263	0.8378	0.7941
	ProtAugment	<b>0.9412</b>	<b>0.9183</b>	<b>0.9412</b>	<b>0.9412</b>	<b>0.9412</b>
	Synonym substitution	0.8235	0.7548	0.8235	0.8235	0.8235
	Tweet style	<b>0.9412</b>	<b>0.9183</b>	<b>0.9412</b>	<b>0.9412</b>	<b>0.9412</b>
<b>ModernBERT</b>	AnnoMI (Base Data)	0.9118	0.8125	0.9040	0.9209	0.9118
	Backtranslation	<b>0.9412</b>	<b>0.9183</b>	<b>0.9412</b>	<b>0.9412</b>	<b>0.9412</b>
	Basic style	0.7941	0.6490	0.7760	0.7740	0.7941
	ProtAugment	0.8235	0.6683	0.8004	0.8137	0.8235
	Synonym substitution	0.8529	0.7308	0.8400	0.8475	0.8529
	Tweet style	0.9118	0.8558	0.9097	0.9097	0.9118

Table 3: Classification Metrics for different models and NLP-augmented data

Model	Semantic Drift			Hallucination			Session Length			Prompt Compliance			Fluency & Emotion		
	High	Avg	Low	High	Avg	Low	Longer	Equal	Shorter	High	Avg	Low	High	Avg	Low
ChatGPT			✓			✓		✓		✓			✓		
DeepSeek			✓			✓			✓	✓				✓	
Gemini		✓			✓				✓		✓		✓		
Llama		✓			✓				✓		✓			✓	
Mistral Nemo		✓			✓				✓		✓			✓	
Nemotron Mini			✓			✓	✓			✓				✓	
Phi-4		✓			✓				✓		✓		✓		
Qwen		✓				✓			✓		✓			✓	
SmolLM2		✓			✓		✓				✓			✓	

Table 4: Evaluation scheme for analyzing the quality of synthetic MI dialogues generated for each LLM.

<b>Model</b>	<b>Semantic Drift</b>	<b>Hallucination</b>	<b>MI Session Length</b>	<b>MI Adherence</b>	<b>Fluency &amp; Emotion</b>
<b>ChatGPT</b>	Very low drift, preserves original intent	Minimal hallucination, rarely adds content	Matches original length closely	Excellent compliance with MI structure	Highly fluent with strong emotional understanding
<b>DeepSeek</b>	Moderate drift, restructures reflections	Rare hallucination but simplifies responses	Slightly condensed compared to original	Adheres to MI but compresses content	Fluent but lacks natural expressiveness
<b>Gemini Flash</b>	Minimal drift, sometimes reinterprets tone	Rare, but adds motivational cues	Slightly longer due to clarifications	Strong adherence but slightly directive	Highly fluent with emotional depth
<b>Llama</b>	Moderate drift, explains rather than reflects	Occasionally adds interpretations	Inconsistent, some responses very brief	Follows MI structure but shifts therapist tone	Fluent but robotic in phrasing
<b>Mistral Nemo</b>	Moderate drift, rewords client intent	Occasionally adds emotional details	Inconsistent, some responses long	Adheres to MI but can be directive	Fluent but lacks emotional nuance
<b>Nemotron Mini</b>	Low-Moderate drift, moderate rewording	Rare hallucination but adds details	Mostly longer than original	Good adherence but over explanation	Fluent but lacks emotional context
<b>Phi4</b>	Low-moderate drift, reframes statements	Adds clarifications, making responses longer	Slightly longer than original	Strong MI adherence but overuses affirmations	Well-structured with moderate emotional depth
<b>Qwen2.5</b>	Low drift, over-summarizes at times	Avoids hallucination but loses nuance	Shorter, condenses responses	Good adherence but loses open-ended depth	Fluent but less expressive
<b>Qwen2.5</b>	Minimal drift, slight rewording	Doesn't add but removes details	Generally shorter than original	Good compliance but oversimplifies	Fluent but lacks emotional depth
<b>Smollm2</b>	Moderate-High drift, restructures statements	Moderate hallucination, adds extra details	Mostly longer than original	Adheres to MI but expands content	Fluent but lacks emotional context

Table 5: An intuitive evaluation of LLMs based on the LLM-augmented data performance.