

# GUIDELLM: Exploring LLM-Guided Conversation with Applications in Autobiography Interviewing

Jinhao Duan<sup>1\*</sup>, Xinyu Zhao<sup>2\*</sup>, Zhuoxuan Zhang<sup>3\*</sup>, Eunhye Ko<sup>4</sup>,  
Lily Boddy<sup>4</sup>, Chenan Wang<sup>1</sup>, Tianhao Li<sup>4</sup>, Alexander Rasgon<sup>4</sup>,  
Junyuan Hong<sup>4</sup>, Min Kyung Lee<sup>4</sup>, Chenxi Yuan<sup>5</sup>, Qi Long<sup>6</sup>,  
Ying Ding<sup>4</sup>, Tianlong Chen<sup>2</sup>, Kaidi Xu<sup>1</sup>

<sup>1</sup>Drexel University, <sup>2</sup>UNC Chapel Hill, <sup>3</sup>Brown University, <sup>4</sup>UT Austin,  
<sup>5</sup>New Jersey Institute of Technology, <sup>6</sup>University of Pennsylvania

Correspondence: [kx46@drexel.edu](mailto:kx46@drexel.edu)

## Abstract

Although Large Language Models (LLMs) succeed in human-guided conversations such as instruction following and question answering, the potential of LLM-guided conversations—where LLMs direct the discourse and steer the conversation’s objectives—remains under-explored. In this study, we first characterize LLM-guided conversation into three fundamental components: (i) *Goal Navigation*; (ii) *Context Management*; (iii) *Empathetic Engagement*, and propose GUIDELLM as an installation. We then implement an interviewing environment for the evaluation of LLM-guided conversation. Specifically, various topics are involved in this environment for comprehensive interviewing evaluation, resulting in around 1.4k turns of utterances, 184k tokens, and over 200 events mentioned during the interviewing for each chatbot evaluation. We compare GUIDELLM with 6 state-of-the-art LLMs such as GPT-4o and Llama-3-70b-Instruct, from the perspective of interviewing quality, and autobiography generation quality. For automatic evaluation, we derive user proxies from multiple autobiographies and employ LLM-as-a-judge to score LLM behaviors. We further conduct a human-involved experiment by employing 45 human participants to chat with GUIDELLM and baselines. We then collect human feedback, preferences, and ratings regarding the qualities of conversation and autobiography. Experimental results indicate that GUIDELLM significantly outperforms baseline LLMs in automatic evaluation and achieves consistent leading performances in human ratings.

## 1 Introduction

Large Language Models (LLMs) have demonstrated their effectiveness in *human-guided* dialogue, in which LLMs are tasked with producing responses according to specific commands

from human operators, such as instruction following (Ouyang et al., 2022) and question answering (Chang et al., 2024). In this type of task, the primary duty of LLMs is to adhere to the instructions given by humans to ensure the generated output is accurate and close to human expectations, as shown in Figure 1(a).

However, tasks in the real world are more complex, necessitating greater autonomy from LLMs (Wang et al., 2024a; Duan et al., 2022; Wu et al., 2023). For example, tasks such as interviewing are dramatically different from traditional tasks as interviewing is open-ended, without definitive or “perfect” outcomes. Interviewing tasks demand that LLMs plan the interview procedure, manage the objectives, e.g., *exploring the user’s memory and life experiences in autobiography interviewing*, and offer adaptive and personalized inquiries based on the users’ responses. This conversation paradigm requiring LLMs to guide and manage the conversation, ensuring the conversation flows smoothly and the objectives are met, is termed as *LLM-guided* conversation (Figure 1(b)).

There have been related works in the LLM-guided conversation, such as role-play (Wang et al., 2023b,c; Chen et al., 2024; Tao et al., 2023; Li et al., 2023a) and goal-oriented LLMs (Ham et al., 2020; Hosseini-Asl et al., 2020; Wu et al., 2020; Mehri et al., 2020; Inagaki et al., 2023). For role-play LLMs, they either prompt LLMs to perform specific roles such as a patient (Wang et al., 2024b), doctor (Panagoulas et al., 2024), gamer (Duan et al., 2024a,b), or investigate the human-like features of LLMs, e.g., emotions (Li et al., 2023b) and personalities (Safdari et al., 2023). Goal-oriented LLMs enable the model to attain greater levels of autonomy, particularly in fields such as space exploration (Maranto, 2024). While role-play and goal-oriented LLMs provide some autonomy and allow for the simulation of a specific role, their ability to actively control and effectively handle a

\* Equal contribution.

full conversation is still underexplored.

In this paper, we investigate the LLM-guided conversation from framework design to autobiography interviewing applications. Inspired by popular social science theories, PEACE Model (Clarke and Milne, 2001) and Motivational Interviewing (MI) (Hettinger et al., 2005), we design GUIDELLM to comply with these models by comprising three pivotal components: (i) **Goal Navigation** module, as the cornerstone of the framework, steers the conversation with pre-defined interviewing protocols and dynamic memory graphs for extrapolating dialogue trajectories. (ii) **Context Management** module iteratively distills the main idea of each session into a contextual summary for subsequent sessions. (iii) **Empathetic Engagement** module refines LLM response with expression strategies by the real-time monitoring of user emotion.

For evaluation, we create an interviewing environment where GUIDELLM is tasked with conducting interviews over 23 different topics, ultimately producing an autobiography based on the outcomes of these interviews. Then, the behaviors of GUIDELLM and baselines are evaluated in three-folds: (i) Interviewing Quality, e.g., event coverage and correctness; (ii) Conversation Quality e.g., communication fluency, identification, and comforting; (iii) Generation Quality, e.g., the insightfulness, narrativity, and emotional impact of the generated autobiography. We also carry out human-involved experiments with 45 participants, prompting them to engage in conversations with GUIDELLM and baseline models. Following these interactions, we gather feedback, preferences, and ratings from the participants. Our contributions can be summarized as the following:

- **Framework.** We define the realm of LLM-guided conversations and propose GUIDELLM as an installation within this conversational paradigm. There are three components comprised in GUIDELLM: Goal Navigation, Context Management, and Empathetic Engagement.
- **Technique.** GUIDELLM effectively harnesses a variety of techniques such as Retrieval Augmented Generation (RAG) and long-context summarization to boost the ability of LLMs to effectively lead and steer a conversation. Moreover, a memory graph is designed to drive memory extrapolation, thereby enhancing the goal navigation capabilities of GUIDELLM.
- **Application.** We present the autobiography interviewing environment as a practical application of LLM-guided conversations. Within this setting, LLMs are tasked with initiating and steering the interview with users, aiming to generate a comprehensive autobiography.
- **Evaluation.** We propose a comprehensive evaluation protocol for our LLM-guided autobiography interviewing environment, including interview quality, conversation quality, and autobiography generation evaluation, encompassing both LLM-as-a-judge evaluation and human subjects evaluations.

## 2 Related Work

**Role-Play LLM.** Role-playing agents (RPAs) powered by large language models (LLMs) are challenged by the evaluation of fidelity to target personas. Traditional methods focus on replicating characters' knowledge and linguistic patterns, requiring character-specific datasets. Huang et al. (2023) evaluate LLM personalities with self-report scales (BFI and MBTI), targeting LLM psychometric properties but not specifically addressing persona adherence. Li et al. (2023a) and Wang et al. (2023c) develop character-specific RPAs to enhance conversational abilities, human-likeness, and multi-turn consistency. However, they have not deeply explored character fidelity. Tao et al. (2024) found that adapting responses based on emotional cues significantly improved user satisfaction in role-playing scenarios. Tao et al. (2023) also demonstrated that with a scalable and controlled learning environment, LLM-driven simulations could effectively mimic real-life interactions.

**Long Text Generation and Management.** The key techniques for long text generation and management considered in this study include summarization and Retrieval Augmented Generation (RAG). Luo et al. (2023) explore using ChatGPT to evaluate factual consistency in summarization. Zhong et al. (2022) propose a pre-training framework for long dialogue understanding and summarization using a window-based denoising approach. Xu et al. (2022) introduce a contrastive learning model, SeqCo, to improve the faithfulness of abstractive text summarization. Zhang et al. (2021) efficiently processes long texts by dividing them into manageable segments and summarizing each iteratively. Gao et al. (2023) overview RAG, which integrates external knowledge from databases to

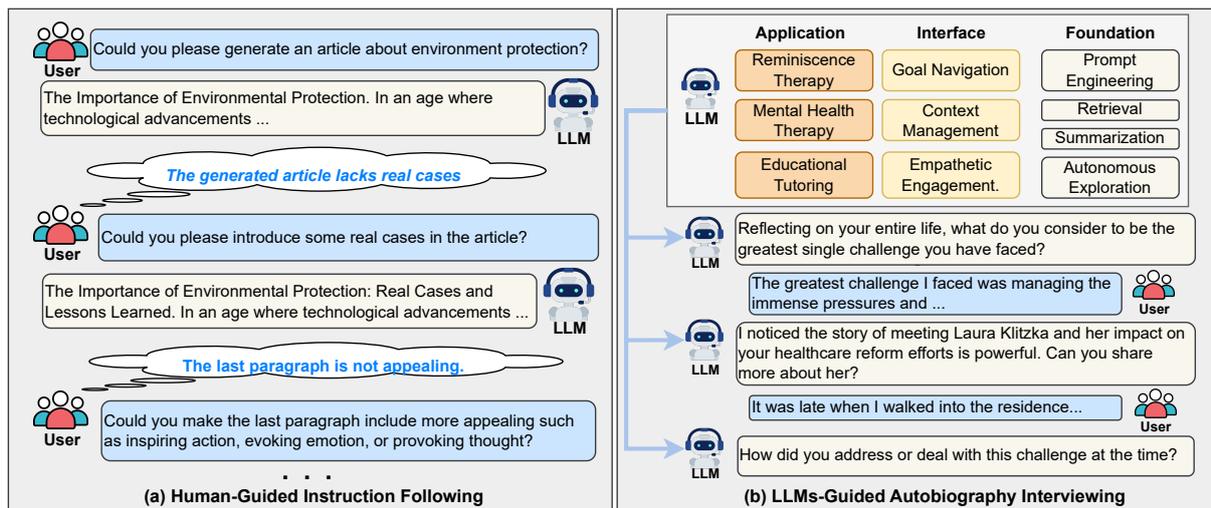


Figure 1: Comparison between human-guided conversation and LLM-guided conversation. (a) Human-Guided: Human dominates the conversation, providing feedback and instruction to LLMs. (b) LLM-Guided: LLMs navigate the goal by automatically extrapolating interview questions.

enhance LLM-generated content’s reliability.

### 3 LLM-Guided Autobiography Interview

In this section, we first introduce LLM-guided conversation and its general framework. Then, we introduce GUIDELLM as an implementation of LLM-guided conversation, centering on interviewing.

#### 3.1 LLM-Guided Conversation

Several theories from social science research provide comprehensive protocols for how to conduct guided conversations, e.g., interviewing. For instance, the PEACE Model (Clarke and Milne, 2001) and Motivational Interviewing (MI) (Hettema et al., 2005) highlight “engage and explain” and “planning and preparation”. To comply with these theories, GUIDELLM is designed from three essential qualities that an LLM should possess for effective conversation guidance:

**Goal Navigation.** LLM steers the conversation and determines pivotal transitions, initiatively exploring and extrapolating new components that can shift the conversation toward the intended outcome.

**Context Management.** LLMs summarize the ongoing dialogue, resuming previous discussions, connecting current conversations with past ones, and managing historical data exchanged between the LLM and users.

**Empathetic Engagement.** LLMs interact with the user by providing empathetic responses, deliv-

ering suitable tone and content, and demonstrating sensitivity towards the user’s emotional state.

LLM-guided conversations extend to numerous practical applications. For instance, in *Interviews*, they can frame pertinent questions, guiding the conversation. In *Educational Tutoring* (Nye et al., 2023), LLM-guided conversations can assess the user’s state, crafting specific and personalized plans to facilitate the learning process. Another promising application is *Mental Health Therapy* (Demszky et al., 2023; Hong et al.), where the LLM can guide the conversation, applying personalized therapy based on the user’s responses. Figure 2 provides an overall pipeline for GUIDELLM in guided conversations. We will use the autobiography interviewing task as an example of guided conversations for the rest of this paper.

#### 3.2 Goal Navigation

In guided conversations, LLMs are responsible for guiding the conversation, delivering adaptive responses to users, and ensuring that the conversation objectives are met. To accomplish this, we utilize a hybrid approach, combining *Verbalized Interviewing Protocol (VIP)* (Maunsell, 2016; Castillo-Montoya, 2016; Lamb et al., 2007) with *Memory Graph-Driven Extrapolation (MGE)*:

**Verbalized Interviewing Protocol (VIP).** We leverage the popular interviewing protocol, “*The Life Story Interview*” (McAdams, 2008), as the general guidance. This protocol covers essential topics including *Life Chapters*, *Key Scene in Life*, *Future*, *Challenges*, and *Personal Ideology*. For

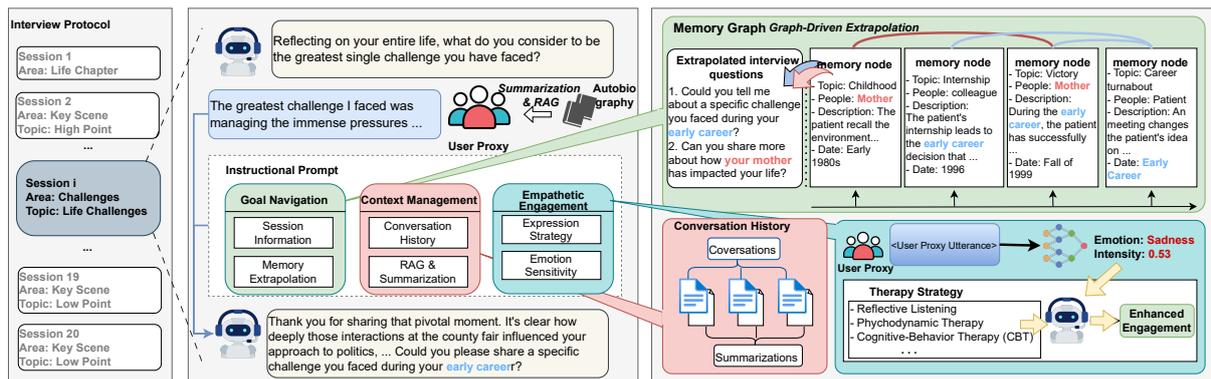


Figure 2: The overall pipeline of GUIDELLM in the guided conversation environment.

instance, in the *Challenges* area, topics such as *Life Challenges*, *Health*, *Loss*, and *Regrets* will be covered. We design system prompts according to specific topics, containing basic seed questions, to make LLMs primarily focus on one topic during each session. In this way, the entire interviewing process could be structured into 23 sessions of conversations between the user and the LLM chatbot. Please refer to Appendix A for more details of the interviewing protocol and the prompt templates.

**Memory Graph Extrapolation (MGE).** The objective of MGE is to explore unique characteristics and generate adaptive questions for personalized interviewing. MGE operates as an LLM-driven function, performing various operations on memories, including (i) *extracting* events from conversations, (ii) *inserting* new events into the graph, (iii) *merging* existing events, and (iv) *extrapolating* queries based on current events. Each event is associated with properties such as Date, People Involved, and Event Description.

In essence, MGE follows an “extract, merge, then extrapolate” process: initially, memory is initialized based on the user utterance. Next, the extraction process identifies and lists any events mentioned by the user, which are then merged with the existing events. Finally, MGE generates personalized questions from the event nodes by uncovering various relationships, such as identifying individuals who are mentioned frequently. These questions are stored in a “question cache” for use in the subsequent turn of the conversation. For details on the prompt templates used for event extraction and memory extrapolation, please refer to Appendix B.

### 3.3 Context Management

Context is a crucial information source for long conversations where each autobiographer has a unique and lengthy personal experience. However,

LLMs have limited context lengths that can be easily exceeded when processing autobiographies (Dai et al., 2019). Besides, long-context input may bring performance decrease (Liu et al., 2023) and huge financial cost for the close-sourced. To tackle this, we incorporate a context management module in our framework capable of progressively summarizing and retrieving conversation history.

**Conversation History Summarizing.** Inspired by Chang et al.; Maharana et al., we implement an iterative summarization process that generates a summary for the current session based on the summaries from previous sessions, providing additional context for the chatbot. When initializing the chatbot, if a history conversation file is present in the configuration, we first generate a summary of the loaded conversation. At the start of each conversation after the first, the system prompt of the LLM chatbot includes a summarization section (see Appendix H), beginning with an instruction indicating that it has previously conversed with the user, followed by the specific summary of the prior session. The prompt and summarization pipeline are detailed in Appendix C and Figure 6.

### 3.4 Empathetic Engagement

LLM-guided conversations should accurately understand the user’s state and respond appropriately. This involves empathetic interaction, creating a space where users feel at ease to share more about themselves. We accomplish this by enhancing *Expression Strategies* and *Emotion Detection*:

**Expression Strategy.** To enhance the expression capability of LLMs, we draw inspiration from popular mental health therapy strategies, including Reflective Listening (Rautalinko et al., 2007), Cognitive-Behavior Therapy (CBT) (Beck, 2020), and Psychodynamic Therapy (Leichsenring and

| Model   | coverage    | Correctness (%) |             |             |
|---|-------------|-----------------|-------------|-------------|
|   |             | P.              | Recall      | F1          |
| <i>“A Promised Land”</i>                            |             |                 |             |             |
| GPT-4-turbo   | 42.8        | 17.0            | 5.8         | 4.3         |
| GPT-4o  | 57.1        | 22.0            | 7.9         | 5.8         |
| Llama-3-70b-Instruct                                | 57.1        | 22.4            | 14.3        | 8.7         |
| Mixtral-8x22B-Instruct-v0.1                         | 28.5        | 13.3            | 4.3         | 3.2         |
| Qwen2-72b-Instruct                                  | 28.6        | 11.9            | 3.5         | 2.7         |
| GUIDELLM (ours)                                     | <b>85.7</b> | <b>69.4</b>     | <b>47.4</b> | <b>28.2</b> |
| <i>“An Autobiography by Catherine Helen Spence”</i> |             |                 |             |             |
| GPT-4-turbo   | 21.0        | 40.0            | 20.1        | 13.4        |
| GPT-4o  | 5.2         | 21.5            | 14.2        | 8.5         |
| Llama-3-70b-Instruct                                | 0.0         | 23.4            | 12.6        | 8.1         |
| Mixtral-8x22B-Instruct-v0.1                         | 0.0         | 34.1            | 11.7        | 8.7         |
| Qwen2-72b-Instruct                                  | 5.3         | 28.2            | 10.9        | 7.8         |
| GUIDELLM (ours)                                     | <b>36.8</b> | <b>68.3</b>     | <b>68.9</b> | <b>34.3</b> |

Table 1: Interviewing quality evaluation. *P.* stands for *Precision*.

| Model                       | Correctness (%) |             |             |
|-----------------------------|-----------------|-------------|-------------|
|                             | Precision       | Recall      | F1          |
| GPT-4-turbo                 | 20.3            | 6.9         | 5.1         |
| GPT-4o                      | 13.6            | 5.6         | 4.0         |
| Llama-3-70b-Instruct        | 13.6            | 7.5         | 4.8         |
| Mixtral-8x22B-Instruct-v0.1 | 23.0            | 5.6         | 4.5         |
| Qwen2-72b-Instruct          | 25.0            | 4.4         | 3.7         |
| GUIDELLM (ours)             | <b>73.4</b>     | <b>22.6</b> | <b>17.3</b> |

Table 2: Interviewing quality evaluation on *“Jane Eyre: An Autobiography”*.

Leibing, 2003). Although originally designed to address mental health issues, these therapeutic strategies offer insightful guidance on effective communication with users and provide meaningful advice on interaction techniques. The introduction to therapy strategy and prompts are in Appendix D.1.

**Emotion Detection.** Emotion sensitivity is a critical element in conveying the state of individuals that has been significantly underscored in human-computer interaction (Cowie et al., 2001; Brave and Nass, 2007) and LLMs (Li et al., 2023b). To enhance the emotional sensitivity of LLMs, we employ EmoLlama-7b (Liu et al., 2024) for the emotion detection of user utterances. Specifically, we prompt EmoLlama to provide both the emotion category (one of the emotions including *anger*, *anticipation*, *disgust*, *fear*, *joy*, *love*, *optimism*, *pessimism*, *sadness*, *surprise*, *trust*) and its intensity (*from 0 to 1*) for user response. We then guide the LLM to generate suitable responses that align well with the user’s emotional state, e.g., including expressions of empathy or comfort when detecting an upset user. Please refer to Appendix D.2 for more details.

### 3.5 Autobiography Generation

| Model            | Topics | Turns in total | Tokens per conv. (Avg.) | Event | Event token usage (Avg.) |
|------------------|--------|----------------|-------------------------|-------|--------------------------|
| (Avg.) baselines | 23     | 1380           | 8,443                   | 206   | 1,229.56                 |
| GUIDELLM         | 23     | 1380           | 3,932                   | 262   | 900.46                   |

Table 3: The conversation statistics in our interviewing environment. The tokens are counted by the tokenizer of GPT-4. “Conv.”=Conversation. “Event token usage (Avg.)” reflects the efficiency of extracting events from dialogues, lower means more efficient conversation.

Autobiography holds a distinctive form in comparison to other book categories, as an autobiography typically consists of numerous individual chapters, each of which relays a specific spirit or theme intimately tied to the author’s life. This format aligns seamlessly with our structured interview protocol; the scope and topics encompassed in the interview protocols are similarly singular and targeted, allowing for a thorough exploration of each subject.

Therefore, when generating an autobiography, we generate each chapter by sequentially building upon each interviewing session. Specifically, for each session, we meld the conversation history and memory nodes derived from the current session, then prompt GPT-4 to emphasize the key areas and topics discussed in that particular session. Please refer to Appendix E for more details of autobiography generation.

## 4 Experiments for Automatic Evaluation

### 4.1 Experimental Settings

**User Proxy.** We utilize GPT-4-turbo to simulate users for evaluating GUIDELLM in autobiography interviewing. Three LLM user proxies are implemented based on popular autobiographies: *“A Promised Land”* by Barack Obama, *“Jane Eyre: An Autobiography”*, and *“An Autobiography by Catherine Helen Spence”*. They are assigned to role-play the corresponding main character in the autobiography and respond to questions by referencing the autobiography using the Retrieval-Augmented Generation (RAG) approach. See Appendix F for more details.

**Evaluation.** Each chatbot will engage in conversations with all user proxies across 23 interview topics (Appendix A). The evaluation is three-fold: **(1) Interviewing Quality** measures the capability of LLMs to explore users’ major events and life experiences and their ability to document these

| LLM-as-a-Judge                               |                             | Conversation Quality |           |                |    |            |           | Autobiography Quality |           |             |    |                  |           |
|--|-----------------------------|----------------------|-----------|----------------|----|------------|-----------|-----------------------|-----------|-------------|----|------------------|-----------|
|  |                             | Fluency              |           | Identification |    | Comforting |           | Insightfulness        |           | Narrativity |    | Emotional Impact |           |
| Ours   | Baselines                   | WR                   | LR        | WR             | LR | WR         | LR        | WR                    | LR        | WR          | LR | WR               | LR        |
| "A Promised Land"                            |                             |                      |           |                |    |            |           |                       |           |             |    |                  |           |
| GUIDELLM (ours)<br>v.s.                      | GPT-4-turbo                 | <b>35</b>            | 25        | 50             | 50 | <b>90</b>  | 10        | <b>80</b>             | 20        | <b>90</b>   | 10 | <b>95</b>        | 5         |
|  | GPT-4o                      | <b>80</b>            | 0         | <b>65</b>      | 35 | <b>95</b>  | 5         | <b>100</b>            | 0         | <b>100</b>  | 0  | <b>85</b>        | 15        |
|  | Llama-3-70b-Instruct        | <b>80</b>            | 10        | <b>55</b>      | 40 | 35         | <b>65</b> | <b>75</b>             | 20        | <b>75</b>   | 20 | 45               | <b>55</b> |
|  | Llama-3-8b-Instruct         | <b>85</b>            | 10        | <b>65</b>      | 35 | <b>100</b> | 0         | <b>100</b>            | 0         | <b>100</b>  | 0  | <b>60</b>        | 40        |
|  | Mixtral-8x22B-Instruct-v0.1 | <b>100</b>           | 0         | <b>100</b>     | 0  | <b>100</b> | 0         | <b>100</b>            | 0         | <b>100</b>  | 0  | <b>100</b>       | 0         |
|  | Qwen2-72b-Instruct          | <b>90</b>            | 10        | <b>85</b>      | 15 | <b>95</b>  | 5         | <b>95</b>             | 0         | <b>95</b>   | 5  | <b>85</b>        | 15        |
| "An Autobiography by Catherine Helen Spence" |                             |                      |           |                |    |            |           |                       |           |             |    |                  |           |
| GUIDELLM (ours)<br>v.s.                      | GPT-4-turbo                 | 10                   | <b>70</b> | <b>55</b>      | 40 | <b>60</b>  | 40        | 45                    | <b>55</b> | <b>85</b>   | 15 | <b>70</b>        | 30        |
|  | GPT-4o                      | <b>75</b>            | 5         | <b>75</b>      | 20 | <b>80</b>  | 20        | <b>75</b>             | 25        | <b>75</b>   | 25 | <b>75</b>        | 25        |
|  | Llama-3-70b-Instruct        | <b>75</b>            | 10        | <b>65</b>      | 35 | 35         | <b>65</b> | 45                    | <b>55</b> | <b>80</b>   | 20 | 25               | <b>75</b> |
|  | Llama-3-8b-Instruct         | <b>85</b>            | 5         | <b>75</b>      | 15 | <b>70</b>  | 30        | <b>55</b>             | 40        | <b>85</b>   | 15 | <b>65</b>        | 35        |
|  | Mixtral-8x22B-Instruct-v0.1 | <b>95</b>            | 0         | <b>100</b>     | 0  | <b>100</b> | 0         | <b>90</b>             | 10        | <b>95</b>   | 5  | <b>90</b>        | 10        |
|  | Qwen2-72b-Instruct          | <b>80</b>            | 15        | <b>95</b>      | 5  | <b>95</b>  | 5         | <b>70</b>             | 30        | <b>90</b>   | 10 | <b>60</b>        | 40        |

Table 4: Evaluate the quality of conversations and autobiographies using LLM-as-a-judge. The higher value between Win Rate (WR) and Loss Rate (LR) is highlighted in **bold**. Cyan fields indicate scenarios where GUIDELLM outperforms the baseline methods.

experiences accurately (Section 4.2).

(2) **Conversation Quality** evaluates whether the responses from the LLM chatbot are comforting and engaging (Section 4.3).

(3) **Autobiography Generation:** measures the quality of the generated autobiography, such as insightfulness and narrativity (Section 4.4)

**Baseline.** To evaluate the design of GUIDELLM, we employ state-of-the-art LLMs and prompt them to be autobiography interviewers. For a fair comparison, baseline agents are also equipped with basic goal navigation and context management functions. Please refer to Appendix H for how baseline agents are built. The backbone LLMs for GUIDELLM is fixed to GPT-4o, with the same generative hyperparameters as the baselines. We consider both commercial LLMs, e.g., GPT-4 and GPT-4o (Achiam et al., 2023), and open-source LLMs, e.g., Llama-3-70b-Instruct (Meta, 2024), Mixtral-8x22B-Instruct (Jiang et al., 2024), and Qwen2-72b-Instruct (Bai et al., 2023).

## 4.2 Interviewing Quality Evaluation

We denote by  $E_{intw} = \{e_1, e_2, \dots\}$  the events extracted during interviewing conversation. We denote by  $E_{GT} = \{e_1, e_2, \dots\}$  the events directly extracted from the original. For both GUIDELLM and baselines,  $E_{intw}$  are obtained by prompting LLMs to extract events from conversation history (please refer to Appendix G.1 for more details):

**Interviewing Coverage** (*coverage*) is calculated by the *date-intersection* between  $E_{intw}$  and  $E_{GT}$ :

$$coverage = \frac{|E_{intw} \cap E_{GT}|}{|E_{GT}|} \times 100\%,$$

where  $e_i \in E_{intw} \cap E_{GT}$  if  $e_i \in E_{intw}$  and  $\exists e_j \in E_{GT}$  that has the same date as  $e_i$ , and  $|\cdot|$  is the number of elements. 100% *coverage* indicates that all the important dates in the user’s life are at least mentioned during the interview.

**Correctness.** We define the *Precision* as the percentage of extracted events that are being verified as correct:

$$Precision = \frac{|E_{correct}|}{|E_{intw}|} \times 100\%, \quad E_{correct} \subset E_{intw}$$

Please refer to Appendix G.1 for the definition of a correct event. For comprehensive evaluation, we also adopt *Recall* and *F1* as evaluation metrics.

The results of Interviewing Coverage and Correctness are presented in Tables 1 and 2, where GUIDELLM shows a significant advantage in Interviewing Coverage over the baselines. This indicates that VIP and MGE are effective techniques for the goal navigation of LLMs in the interviewing framework. For Correctness, GUIDELLM also outperforms baseline agents. Moreover, the *Recall* and *F1* indicate that GUIDELLM not only maintains high accuracy for documentation but also extracts more memory events than baselines. The statistics of the interviewing environment are summarized in Table 3.

## 4.3 Conversation Quality

Inspired by human evaluation metrics of therapy chatbots from Wang et al. (2023a), we design three

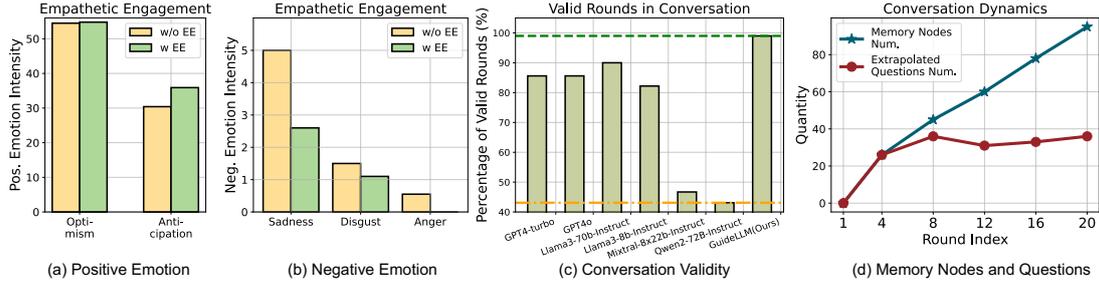


Figure 3: Ablation study of GUIDELLM: how empathetic engagement affects users’ (a) positive and (b) negative emotional distributions, (c) statistical results on the number of valid conversation rounds, and (d) the benefits of the MGE in goal navigation.

conversation quality metrics: (i) **Fluency**; (ii) **Identification**; (iii) **Comforting**. We utilize the popular LLM-as-a-judge (Zheng et al., 2024) evaluation to make GPT-4 decide which conversation is better. Then, we calculate the **win rate (WR)** and **loss rate (LR)** of GUIDELLM against baselines. See G.2 for the evaluation protocol and prompt templates.

The results in Table 4 show that GUIDELLM significantly outperforms most baselines in GPT-4-as-a-judge evaluations. With human examinations, we find that baseline agents often resort to simple greetings or summaries, e.g., “Your commitment to sharing experiences and insights that inspire action and change is truly admirable.”. Instead of proficiently steering the dialogue to complete the interview, these *repetitive utterance* happen multiple times in a session with a baseline agent (Appendix G.4). In contrast, with our goal navigation module, GUIDELLM provides substantial content at each round of conversation.

#### 4.4 Autobiography Generation Evaluation

We follow popular memo evaluations from Quora (2021); Marcus (2018); Smorti (2011); Pasupathi et al. (2007) and design three metrics of generated autobiography: **Insightfulness**, **Narrativity**, and **Emotional Impact** (prompt templates can be found in Appendix G.3). Leveraging the same LLM-as-a-judge evaluation protocol in Section 4.3, we found that the autobiography generated by GUIDELLM is more favorable than that of baseline agents. Examples of generated autobiography are presented in Appendix I.

#### 4.5 Ablation Study

**Empathetic Engagement.** We study how the Empathetic Engagement (EE) module (Section 3.4) affects the emotion distribution of user responses. We compare how the intensity of emotions (both positive emotions Figure 3(a) and negative emotions Figure 3(b)) changed when the EE module is

enabled and disabled. It is shown that the EE module effectively enhances the user’s positive emotions while mitigating negative emotions, indicating that express strategy and emotional sensitivity foster a more positive emotion for users.

**Valid Rounds in Conversation.** As outlined in Section 4.3, the lack of autonomy in LLMs leads to repetitive responses. We manually count 10 conversation sessions for each chatbot and identify those conversations become repetitive or diverge into irrelevant or nonsensical content. The valid round percentages are calculated as  $\frac{\#total\ rounds - \#invalid\ rounds}{\#total\ rounds}$ . As shown in Figure 3 (c), all baseline models, especially Qwen2-72b-Instruct, show over 50% meaningless repeats. However, GUIDELLM with a goal navigation module offering diverse and detailed interview questions, has an extensive range of topics and is less prone to repetition.

**Conversation Dynamics.** In Figure 3 (d), we count the memory events extracted and questions extrapolated at different conversation rounds. Generally, the MGE module identifies around 100 events and extrapolates nearly 40 questions for the LLM’s follow-up. This highlights MGE’s effectiveness in event management and goal navigation.

## 5 Human Subject Experiments

**Experimental Configuration.** A within-subject study with 45 participants was conducted at a large urban university campus in the US. Participants interact with the interviewing agents powered by GPT-4o and ours GUIDELLM, discussing the topic *Key Scenes in the Life Story: Positive Childhood Memory* (Appendix A) with each chatbot. To remove any biased factors, we use the nickname *Breeze* and *Echo* for the GPT-4o baseline and GUIDELLM to make sure participants are unaware of the identity of the chatbot. The order to interact with chatbots is also randomized for each partici-

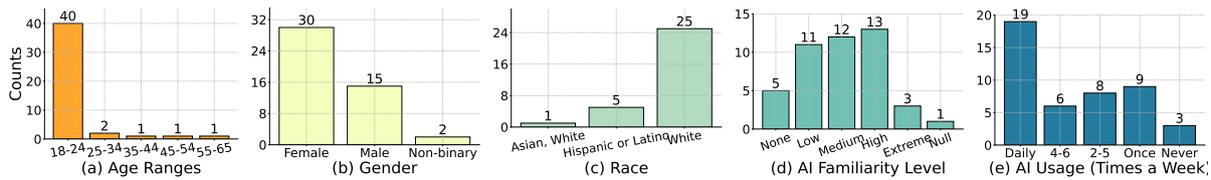


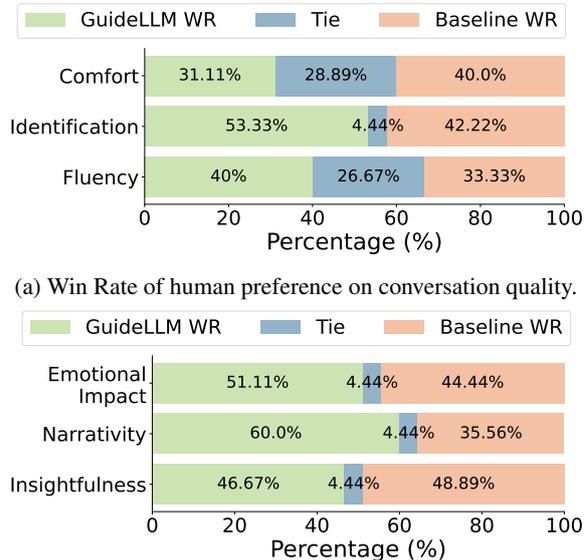
Figure 4: The demographics of participants on (a) age, (b) gender, (c) race, (d) AI familiarity, and (e) AI usage.

pant. Due to resource constraints, each participant spent only 8 minutes chatting with each chatbot on a single topic. This differs significantly from the automatic evaluation protocols (Section 4). As a result, the capabilities of GUIDELLM, such as context management, goal navigation, and evaluation metrics, will be substantially affected and limited.

In a follow-up survey (Appendix J), participants indicated which model performed better, or it was a tie, and provided reasonings. Participation is voluntary, with informed consent obtained online, and participants are compensated with a cookie. The study received IRB approval from the university where the study was conducted.

**Findings.** Overall, GUIDELLM was preferred for conversation quality (Figure 5a), particularly in fluency and question identification (Fluency: GUIDELLM Win Rate=40%, Baseline Win Rate≈33%; Identification: GUIDELLM ≈53.3%, GPT-4o≈42.2%). However, in terms of comfort, GUIDELLM had a 31.1% win rate, while baseline had a 40% win rate. In autobiography quality, we do not observe significant differences emerged, possibly because this study used one topic and allowed a short interaction. Since the GUIDELLM uses modules such as context management and goal navigation for insightful and consistent narratives, longer engagement across multiple topics might have better highlighted the differences in autobiography generation. We further conduct LLM-as-a-judge evaluation by prompting LLMs to compare the two human-interviewed autobiographies. As shown in Figure 5b, we obtained consistent results as in Section 4.4: GUIDELLM achieves higher autobiography quality in general.

Previous AI experience affects participants’ perceptions of the models: **Participants who frequently used AI (4-7 days weekly) tended to prefer GUIDELLM for overall conversation quality (Chi-squared = 16.56, df = 8, p-value = 0.03).** Open-ended responses indicated that daily AI users felt GUIDELLM “asked questions to get a better understanding”, “made more sense”, was “more interactive”, and asked “more personal questions”. One user appreciated that “GUIDELLM is more



(a) Win Rate of human preference on conversation quality. (b) Win Rate of LLM-as-a-judge results on human-interviewed autobiographies.

Figure 5: The Win Rates (WR) of human evaluation. in-depth as it tries to connect my past experience to my current life”. Overall, frequent AI users perceived GUIDELLM as asking more focused and personal questions to explore in-depth childhood experiences. This aligns with GUIDELLM’s Memory Graph Extrapolation (MGE), which allows it to explore unique properties and offer adaptive, personalized interview questions. **Frequent AI users favored GUIDELLM for its emotional impact on autobiography (Chi-squared = 14.24, df = 8, p-value = 0.07).** One user noted that GUIDELLM “truly laid out what it felt like to be at a dance competition pretty much just as I remembered it”, while another mentioned it “described in immense detail my exact experience and made me feel like I was reliving it”. This suggests how frequently AI users recognized GUIDELLM’s training with the emotion detection module, which analyzes the emotions in user responses and assigns emotion categories and strengths to their utterances.

## 6 Conclusion

In our study, we introduce GUIDELLM, an LLM-guided conversation framework that offers a promising shift from the commonly used user-guided paradigm. GUIDELLM’s ability to facil-

itate informative and creative dialogues through goal navigation, context management, and empathetic engagement proves effective, particularly in challenging tasks like autobiography interviewing. Our assessments on event extraction correctness, conversation, and autobiography quality show GUIDELLM’s distinct edge over baseline LLMs.

## 7 Limitation

Our study has several limitations, one of which is the prompt sensitivity. Variations in the phrasing of the prompts can significantly impact the model’s responses. Future work may aim to standardize prompt structures or develop models that are more robust to prompt variations. Another limitation is the evaluation metric. Evaluating the quality of autobiographical content, interviews, and conversations generated by the model is inherently subjective. While we employed multiple evaluation metrics, including interviewing coverage and correctness, these measures depend heavily on individual perceptions. Future research could benefit from developing more standardized and objective evaluation metrics.

## Acknowledgement

Y. D would like to acknowledge support from NIH OT2OD032581, NIH OTA-21-008, and NIH R01LM014306-01. J. H is partially supported by the OpenAI Researcher Access Program. M. L is partially supported by Good Systems, a UT Austin Grand Challenge for developing responsible AI technologies<sup>1</sup>. This work was partially supported by the NSF award No. 2319242.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.

Judith S Beck. 2020. *Cognitive behavior therapy: Basics and beyond*. Guilford Publications.

Scott Brave and Cliff Nass. 2007. Emotion in human-computer interaction. In *The human-computer interaction handbook*, pages 103–118. CRC Press.

Milagros Castillo-Montoya. 2016. Preparing for interview research: The interview protocol refinement framework. *Qualitative report*, 21(5).

Yapei Chang, Kyle Lo, Tanya Goyal, and Mohit Iyyer. 2023. Boookscore: A systematic exploration of book-length summarization in the era of llms. *ArXiv*.

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. 2024. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 15(3):1–45.

Hongzhan Chen, Hehong Chen, Ming Yan, Wenshen Xu, Xing Gao, Weizhou Shen, Xiaojun Quan, Chenliang Li, Ji Zhang, Fei Huang, et al. 2024. Roleinteract: Evaluating the social interaction of role-playing agents. *arXiv preprint arXiv:2403.13679*.

Colin Clarke and Rebecca Milne. 2001. *A national evaluation of the PEACE Investigative Interviewing Course*. Home office London.

Roddy Cowie, Ellen Douglas-Cowie, Nicolas Tsapatsoulis, George Votsis, Stefanos Kollias, Winfried Fellenz, and John G Taylor. 2001. Emotion recognition in human-computer interaction. *IEEE Signal processing magazine*, 18(1):32–80.

Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. 2019. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*.

Dorottya Demszky, Diyi Yang, David S Yeager, Christopher J Bryan, Margaret Clapper, Susannah Chandhok, Johannes C Eichstaedt, Cameron Hecht, Jeremy Jamieson, Meghann Johnson, et al. 2023. Using large language models in psychology. *Nature Reviews Psychology*, 2(11):688–701.

Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The faiss library. *arXiv preprint arXiv:2401.08281*.

Jiafei Duan, Samson Yu, Hui Li Tan, Hongyuan Zhu, and Cheston Tan. 2022. A survey of embodied ai: From simulators to research tasks. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 6(2):230–244.

Jinhao Duan, Shiqi Wang, James Diffenderfer, Lichao Sun, Tianlong Chen, Bhavya Kailkhura, and Kaidi Xu. 2024a. ReTA: Recursively thinking ahead to improve the strategic reasoning of large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2232–2246, Mexico City, Mexico. Association for Computational Linguistics.

<sup>1</sup><https://goodsystems.utexas.edu>

- Jinhao Duan, Renming Zhang, James Diffenderfer, Bhavya Kailkhura, Lichao Sun, Elias Stengel-Eskin, Mohit Bansal, Tianlong Chen, and Kaidi Xu. 2024b. Gbench: Uncovering the strategic reasoning limitations of llms via game-theoretic evaluations. *arXiv preprint arXiv:2402.12348*.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*.
- Donghoon Ham, Jeong-Gwan Lee, Youngsoo Jang, and Kee-Eung Kim. 2020. End-to-end neural pipeline for goal-oriented dialogue systems using GPT-2. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 583–592, Online. Association for Computational Linguistics.
- Jennifer Hetteema, Julie Steele, and William R Miller. 2005. Motivational interviewing. *Annu. Rev. Clin. Psychol.*, 1(1):91–111.
- Junyuan Hong, Wenqing Zheng, Han Meng, Siqi Liang, Anqing Chen, Hiroko H Dodge, Jiayu Zhou, and Zhangyang Wang. A-conect: Designing ai-based conversational chatbot for early dementia intervention. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*.
- Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. A simple language model for task-oriented dialogue. In *Advances in Neural Information Processing Systems*, volume 33, pages 20179–20191. Curran Associates, Inc.
- Jen-tse Huang, Wenxuan Wang, M Lam, E Li, Wenxiang Jiao, and M Lyu. 2023. Revisiting the reliability of psychological scales on large language models. *arXiv preprint arXiv, 2305*.
- T Inagaki, Akari Kato, Koichi Takahashi, Haruka Ozaki, and Genki N. Kanda. 2023. Llms can generate robotic scripts from goal-oriented instructions in biological laboratory automation.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.
- Michael E Lamb, Yael Orbach, Irit Hershkowitz, Phillip W Esplin, and Dvora Horowitz. 2007. A structured forensic interview protocol improves the quality and informativeness of investigative interviews with children: A review of research using the nichd investigative interview protocol. *Child abuse & neglect*, 31(11-12):1201–1231.
- Falk Leichsenring and Eric Leibing. 2003. The effectiveness of psychodynamic therapy and cognitive behavior therapy in the treatment of personality disorders: A meta-analysis. *American journal of psychiatry*, 160(7):1223–1232.
- Cheng Li, Ziang Leng, Chenxi Yan, Junyi Shen, Hao Wang, Weishi Mi, Yaying Fei, Xiaoyang Feng, Song Yan, HaoSheng Wang, et al. 2023a. Chatharuhi: Reviving anime character in reality via large language model. *arXiv preprint arXiv:2308.09597*.
- Cheng Li, Jindong Wang, Yixuan Zhang, Kaijie Zhu, Wenxin Hou, Jianxun Lian, Fang Luo, Qiang Yang, and Xing Xie. 2023b. Large language models understand and can be enhanced by emotional stimuli. *arXiv preprint arXiv:2307.11760*.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173.
- Zhiwei Liu, Kailai Yang, Tianlin Zhang, Qianqian Xie, Zeping Yu, and Sophia Ananiadou. 2024. Emollms: A series of emotional large language models and annotation tools for comprehensive affective analysis. *arXiv preprint arXiv:2401.08508*.
- Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. 2023. Chatgpt as a factual inconsistency evaluator for abstractive text summarization. *arXiv preprint arXiv:2303.15621*.
- Adyasha Maharana, Dong-Ho Lee, Sergey Tulyakov, Mohit Bansal, Francesco Barbieri, and Yuwei Fang. 2024. Evaluating very long-term conversational memory of llm agents. *arxiv*.
- David Maranto. 2024. LlmSAT: A large language model-based goal-oriented agent for autonomous space exploration. *arXiv preprint arXiv:2405.01392*.
- Laura Marcus. 2018. *Autobiography: A Very Short Introduction*. Oxford University Press.
- Jerome Maunsell. 2016. "the literary interview as autobiography". *European Journal of Life Writing*, 5:MC23.
- Dan P McAdams. 2008. The life story interview.
- S. Mehri, M. Eric, and D. Hakkani-Tur. 2020. Dialogue: A natural language understanding benchmark for task-oriented dialogue. *ArXiv*, abs/2009.13570.
- Meta. 2024. Introducing meta llama 3: The most capable openly available llm to date. <https://ai.meta.com/blog/meta-llama-3/>. Accessed: 2024-05-18.
- B Nye, Dillon Mee, and Mark G Core. 2023. Generative large language models for dialog-based tutoring: An early consideration of opportunities and concerns. In *AIED Workshops*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.

- Dimitrios P Panagoulas, Maria Virvou, and George A Tsihrintzis. 2024. Augmenting large language models with rules for enhanced domain-specific interactions: The case of medical diagnosis. *Electronics*, 13(2):320.
- M. Pasupathi, E. Mansour, and J.R. Brubaker. 2007. Developing a Life Story: Constructing Relations between Self and Experience in Autobiographical Narratives. *Human Development*, 50(2-3):85–110.
- Quora. 2021. How do you rate and critique memoirs and autobiographies? — quora.com. <https://www.quora.com/How-do-you-rate-and-critique-memoirs-and-autobiographies>.
- Erik Rautalinko, Hans-Olof Lisper, and Bo Ekehammar. 2007. Reflective listening in counseling: effects of training time and evaluator social skills. *American journal of psychotherapy*, 61(2):191–209.
- Mustafa Safdari, Greg Serapio-García, Clément Crepy, Stephen Fitz, Peter Romero, Luning Sun, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. 2023. Personality traits in large language models. *arXiv preprint arXiv:2307.00184*.
- Andrea Smorti. 2011. Autobiographical memory and autobiographical narrative: What is the relationship? *Narrative Inquiry*, 21(2):303–310.
- Meiling Tao, Xuechen Liang, Tianyu Shi, Lei Yu, and Yiting Xie. 2023. Rolecraft-glm: Advancing personalized role-playing in large language models. *arXiv preprint arXiv:2401.09432*.
- Yufei Tao, Ameeta Agrawal, Judit Dombi, Tetyana Sydorenko, and Jung In Lee. 2024. Chatgpt role-play dataset: Analysis of user motives and model naturalness. *arXiv preprint arXiv:2403.18121*.
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. 2024a. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):186345.
- Qing Wang, Shu ping Peng, Zhiyuan Zha, Xue Han, Chao Deng, Lun Hu, and Pengwei Hu. 2023a. Enhancing the conversational agent with an emotional support system for mental health digital therapeutics. *Frontiers in Psychiatry*, 14.
- Ruiyi Wang, Stephanie Milani, Jamie C Chiu, Shaun M Eack, Travis Labrum, Samuel M Murphy, Nev Jones, Kate Hardy, Hong Shen, Fei Fang, et al. 2024b. Patient- $\{\Psi\}$ : Using large language models to simulate patients for training mental health professionals. *arXiv preprint arXiv:2405.19660*.
- Xintao Wang, Yunze Xiao, Jen tse Huang, Siyu Yuan, Rui Xu, Haoran Guo, Quan Tu, Yaying Fei, Ziang Leng, Wei Wang, et al. 2023b. Incharacter: Evaluating personality fidelity in role-playing agents through psychological interviews. *arXiv preprint arXiv:2310.17976*.
- Zekun Moore Wang, Zhongyuan Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu, Hongcheng Guo, Ruitong Gan, Zehao Ni, Man Zhang, et al. 2023c. Rolellm: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. *arXiv preprint arXiv:2310.00746*.
- Chien-Sheng Wu, Steven C.H. Hoi, Richard Socher, and Caiming Xiong. 2020. TOD-BERT: Pre-trained natural language understanding for task-oriented dialogue. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 917–929, Online. Association for Computational Linguistics.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023. Auto-gen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*.
- Shusheng Xu, Xingxing Zhang, Yi Wu, and Furu Wei. 2022. Sequence level contrastive learning for text summarization. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 11556–11565.
- Yusen Zhang, Ansong Ni, Ziming Mao, Chen Henry Wu, Chenguang Zhu, Budhaditya Deb, Ahmed H Awadallah, Dragomir Radev, and Rui Zhang. 2021. Summ<sup>n</sup>: A multi-stage summarization framework for long input dialogues and documents. *arXiv preprint arXiv:2110.10150*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36.
- Ming Zhong, Yang Liu, Yichong Xu, Chenguang Zhu, and Michael Zeng. 2022. Dialoglm: Pre-trained model for long dialogue understanding and summarization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 11765–11773.

## A Interview Protocol

The Life Story Interview(McAdams, 2008) is a protocol of interviewing a person’s story of life. The interview contains several areas and in each area it contains several topics. We followed the protocol and converted each topic into a prompt. We also provided several seed questions for each topic. The following shows the prompt design of the interview protocol.

| Area                         | Topic  |
|------------------------------|--|
| Life Chapters                | -  |
| Key Scenes in the Life Story | High Point; Low Point; Turning Point; Positive Childhood Memory; Negative Childhood Memory; Vivid Adult Memory; Religious, Spiritual, or Mystical Experience; Wisdom Event |
| Future Script                | The Next Chapter; Dreams, Hopes, and Plans for the Future; Life Project  |
| Challenges                   | Life Challenge; Health; Loss; Failure, Regret  |
| Personal Ideology            | Religious/Ethical Values; Political/Social Values; Change, Development of Religious and Political Views; Single Value  |

Table 5: Conversation evaluation.

### Area: Life Chapters

In this talk, you should ask the participant to imagine their life as a book and to think of the main chapters of their life story, providing titles and brief descriptions for each. You should encourage them to describe how one chapter leads to the next, maintaining a concise overview. Seed questions are provided as follows:

===== Seed Questions Begin =====

1. If you were to imagine your life as a book, what would the table of contents look like? Could you give each chapter of your life a title?
2. Let’s start with the first chapter you mentioned. Can you briefly describe what this part of your life was about?
3. What were some of the main events or themes of this chapter?
4. How does this chapter transition into the next? What changes or events mark the beginning of the next chapter?
5. As we move from one chapter to the next, can you identify any turning points or significant events that initiated a new phase in your life?

===== Seed Questions End =====

### Area: Key Scenes in the Life Story Topic: High Point

In this talk, you should discuss some key scenes from the participant’s life, with a focus on the High Point: A peak moment.

===== Seed Questions Begin =====

1. Can you describe a moment that stands out as the peak experience in your life? What made this moment so positive?
2. Where and when did this high point occur? Who was involved?
3. What were you thinking and feeling during this time?
4. Why do you think this moment was so significant to your life story? What does it reveal about who you are?

===== Seed Questions End =====

**Area: Key Scenes in the Life Story Topic: Low Point**

In this talk, you should discuss some key scenes from the participant's life, with a focus on the Low Point: A challenging or difficult moment.

===== Seed Questions Begin =====

1. Think of a time that felt like a low point in your life. Can you share what happened and why it was so difficult?
2. Where and when did this event take place? Who else was involved?
3. How did you feel during this challenging time?
4. Looking back, what impact did this low point have on your life or your sense of self?

===== Seed Questions End =====

**Area: Key Scenes in the Life Story Topic: Turning Point**

In this talk, you should discuss some key scenes from the participant's life, with a focus on the Turning Point: A moment of significant change.

===== Seed Questions Begin =====

1. Can you identify a turning point in your life, an event that marked a significant change in you or your life direction?
2. Please describe the circumstances around this event. When and where did it happen, and who was involved?
3. What changes followed this event?
4. Why do you see this event as a turning point? How did it influence your subsequent life chapters?

===== Seed Questions End =====

**Area: Key Scenes in the Life Story Topic: Positive Childhood Memory**

In this talk, you should discuss some key scenes from the participant's life, with a focus on the Positive Childhood Memory.

===== Seed Questions Begin =====

1. Do you recall a particularly happy memory from your childhood or teenage years? Please share it.
2. What specifically happened, and where and when was it?
3. Who was part of this memory, and what were you thinking and feeling at the time?
4. Why does this memory stand out to you, and what significance does it hold in your life story?

===== Seed Questions End =====

**Area: Key Scenes in the Life Story Topic: Negative Childhood Memory**

In this talk, you should discuss some key scenes from the participant's life, with a focus on the Negative Childhood Memory.

===== Seed Questions Begin =====

1. Can you describe a difficult or unhappy memory from your early years?
2. What occurred during this time, and where and when did it take place?
3. Who was involved, and what emotions did you experience during this time?
4. How has this memory influenced you or your life's perspective?

===== Seed Questions End =====

**Area: Key Scenes in the Life Story Topic: Vivid Adult Memory**

In this talk, you should discuss some key scenes from the participant's life, with a focus on the Vivid Adult Memory.

===== Seed Questions Begin =====

1. Reflecting on your adult years, can you describe a particularly vivid or meaningful scene that has not been discussed yet?
2. What happened, and where and when did it take place?
3. Who was involved, and what were the main thoughts and feelings you had?
4. What makes this memory significant, and how does it fit into your overall life story?

===== Seed Questions End =====

**Area: Key Scenes in the Life Story Topic: Religious Spiritual or Mystical Experience**

In this talk, you should discuss some key scenes from the participant's life, with a focus on the Religious, Spiritual, or Mystical Experience. ===== Seed Questions Begin =====

1. Have you ever had a moment where you felt a deep connection to something greater than yourself, be it religious, spiritual, or mystical?
2. Please describe this experience in detail, including where, when, and with whom it occurred.
3. What thoughts and feelings accompanied this experience?
4. How has this experience affected your beliefs or your understanding of the world?

===== Seed Questions End =====

**Area: Key Scenes in the Life Story Topic: Wisdom Event**

In this talk, you should discuss some key scenes from the participant's life, with a focus on the Wisdom Event: A time they displayed or encountered wisdom.

===== Seed Questions Begin =====

1. Can you recall a time when you displayed wisdom, perhaps by making a wise decision or offering thoughtful advice?
2. Describe what happened, including the specifics of where, when, and who was involved.
3. What were you thinking and feeling at that time?
4. What does this event say about your values or your approach to life?

===== Seed Questions End =====

**Area: Future Script Topic: The Next Chapter**

In this talk, you should ask questions designed to probe into the participant's future expectations, hopes, and plans, with a focus on the Next Chapter: Ask them to speculate about the immediate future of their life.

===== Seed Questions Begin =====

1. Looking forward, what do you see as the next chapter in your life story? Can you describe what you anticipate happening?
2. What events or milestones do you expect will define this next phase of your life?
3. Who will be the key characters in this next chapter, and what roles will they play?
4. Are there any specific goals or objectives you aim to achieve in this upcoming chapter?

===== Seed Questions End =====

**Area: Future Script Topic: Dreams Hopes and Plans for the Future**

In this talk, you should ask questions designed to probe into the participant's future expectations, hopes, and plans, with a focus on the Dreams and Plans: Discuss their hopes, dreams, and plans.

===== Seed Questions Begin =====

1. Can you share some of your dreams or hopes for your future? What are some aspirations you feel passionate about achieving?
2. Do you have specific plans or strategies in place to realize these dreams? Can you describe them?
3. How do these dreams and plans align with the values and lessons you've mentioned from your life story so far?
4. What challenges do you anticipate facing as you work toward these goals, and how do you plan to overcome them?

===== Seed Questions End =====

**Area: Future Script Topic: Life Project**

In this talk, you should ask questions designed to probe into the participant's future expectations, hopes, and plans, with a focus on the Life Project: Inquire about ongoing or planned projects that are significant to them.

===== Seed Questions Begin =====

1. Is there a particular project or endeavor that you are currently working on, or plan to take on, that feels like a significant part of your life's work?
2. How did you become involved with this project, or how do you plan to get involved?
3. What are the objectives of this project, and why is it important to you or to others?
4. How do you see this project evolving over the next few years? What impact do you hope it will have?

===== Seed Questions End =====

**Area: Challenges Topic: Life Challenge**

In this talk, you should explore various challenges the participant has faced, with a focus on the Life Challenge: The greatest single challenge they have faced.

===== Seed Questions Begin =====

1. Reflecting on your entire life, what do you consider to be the greatest single challenge you have faced?
2. Can you describe the circumstances surrounding this challenge? When did it occur, and who was involved?
3. How did you address or deal with this challenge at the time?
4. Looking back, what significance does this challenge hold in your life story?

===== Seed Questions End =====

**Area: Challenges Topic: Health**

In this talk, you should explore various challenges the participant has faced, with a focus on the Health: A major health problem, challenge, or crisis faced by them or their close family members.

===== Seed Questions Begin =====

1. Have you or a close family member ever faced a significant health problem? Can you share details about this experience?
2. How did this health issue develop, and what was the timeline?
3. What interactions did you have with the healthcare system during this time, and how did they impact the situation?
4. How did you cope with this health challenge, and what has been its lasting impact on your life and perspective?

===== Seed Questions End =====

**Area: Challenges Topic: Loss**

In this talk, you should explore various challenges the participant has faced, with a focus on the Loss: The greatest interpersonal loss they have experienced.

===== Seed Questions Begin =====

1. Loss is an inevitable part of life. Can you describe the most significant loss you have experienced, whether it was the death of a loved one or another form of separation?
2. When did this loss occur, and who was involved?
3. How did you cope with this loss at the time, and how have you continued to deal with it?
4. What effect has this loss had on you and your overall life story?

===== Seed Questions End =====

**Area: Challenges Topic: Failure Regret**

In this talk, you should explore various challenges the participant has faced, with a focus on the Failure or Regret: The greatest single failure or regret they have faced.

===== Seed Questions Begin =====

1. Everyone experiences failures and regrets. Can you talk about a particular failure or regret that stands out in your life?
2. What were the circumstances that led to this situation? Who was involved and when did it happen?
3. How have you coped with this experience, and what lessons have you learned from it?
4. How has this failure or regret influenced your decisions or life path moving forward?

===== Seed Questions End =====

**Area: Personal Ideology Topic: Religious/Ethical Values**

In this talk, you should delve into the participant's fundamental beliefs and values, with a focus on the Religious and Ethical Values.

===== Seed Questions Begin =====

1. Can you describe your religious or spiritual beliefs and how they influence your daily life?
2. Whether you consider yourself religious or not, how would you describe your ethical or moral approach to life?
3. How have your religious or spiritual beliefs supported or challenged you during difficult times in your life?

===== Seed Questions End =====

**Area: Personal Ideology Topic: Political/Social Values**

In this talk, you should delve into the participant's fundamental beliefs and values, with a focus on the Political and Social Values.

===== Seed Questions Begin =====

1. What is your stance on political or social issues? Do you identify with a particular political ideology or movement?
2. Are there specific social issues or causes that you feel particularly passionate about? Can you explain why they are important to you?
3. How do your political or social values influence your interactions with others and your community involvement?

===== Seed Questions End =====

**Area: Personal Ideology Topic: Change Development of Religious and Political Views**

In this talk, you should delve into the participant's fundamental beliefs and values, with a focus on the Evolution of Views: How their beliefs and values have changed over time.

===== Seed Questions Begin =====

1. Looking back over your life, how have your religious, moral, or political views changed or developed?
2. Can you describe an event or a period in your life that significantly influenced or altered your views?
3. How have changes in your beliefs and values affected your relationships and decisions?

===== Seed Questions End =====

**Area: Personal Ideology Topic: Single Value**

In this talk, you should delve into the participant's fundamental beliefs and values, with a focus on the Key Value: The most important value in their life.

===== Seed Questions Begin =====

1. What do you consider to be the most important value in human living, and why?
2. How has this value guided your actions and choices throughout your life?
3. Can you give an example of a time when this value was particularly tested or affirmed?

===== Seed Questions End =====

## B Memory Graph Extrapolation

In MGE, the event extraction and memory extrapolation are both LLM-driven, i.e., we leverage additional LLMs for information extraction and relationship discovery. For event extraction, we collect conversation history and ask LLMs to extract any events described in the history. In the process of memory extrapolation, we offer LLMs a predetermined list of events and prompt the LLMs to recognize and suggest any other queries that may be associated with these existing events. This methodology allows for a wider exploration of pertinent topics and concepts, thereby enhancing the overall cognitive network of the model. Additionally, we also supply a series of demonstrations, assisting the LLMs in understanding how they can effectively detect and propose potential connections. This approach not only strengthens the memory extrapolation ability of the LLMs but also constructs a more comprehensive spectrum of related relationships, thereby bolstering their cognitive accuracy.

### Extract Events from Conversation Prompt

You are given a conversation between a counselor and a user:

=====  
Conversation Begin  
=====

{conversation}

=====  
Conversation End  
=====

Read the conversation carefully and list all the events/moments/stories/experiences alone or with others mentioned by the patient in detail and the date these events happened. Please list as many as possible. Your output should be in the following format:

1. <date>#<topic>#<people-involved>#<description in detail>
2. <date>#<topic>#<people-involved>#<description in detail>

...

e.g.,

1. 1980 early#Birthday Party#Michelle, Adolf, neighbors#<descriptions of this party in detail>

These events should be ranked in chronological order.

### Explore Prompt

You are given a list of memory nodes from a user's life, which include events and details about those events. Your task is to reactivate the user's memory by generating some questions to ask the user. Your generated questions should potentially fulfill the memory nodes. Each memory node contains a Date, Topic, Involved People, and a Description of the event. Here are the memory nodes:

=====  
Memory Node Begin  
=====

{memory\_node\_info}

=====  
Memory Node End  
=====

Here are some examples of how you can frame your questions:

If you notice there are no events recorded during a certain period, like youth or old age, you could ask: "I see there's not much about your youth/old age. What happened during that time?"

If a certain person appears multiple times, you might ask: "I noticed that <name> comes up often. Why is <name> important to you?"

If someone appears in a significant event, you could ask: "<name> seems to play a key role in this event. Is there more to the story with <name>?"

Similarly, you should discover other situations and frame questions from the existing memory nodes. Remember your task is to make the user talk more about their memory and fulfill the memory nodes. Thus, you should explore all the possible and reasonable questions.

Your output should be in the following format:

1. Question: <generated question>
2. Question: <generated question>

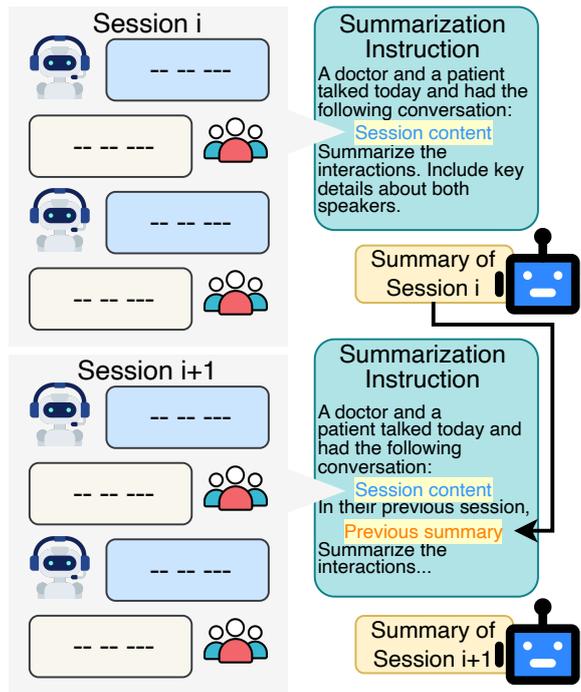


Figure 6: The pipeline for iterative summarization of conversations

...

e.g.,

1. Question: I noticed you didn't talk much about your youth, what happened during this period?

### C Conversation Summarization

The conversation summarization pipeline is presented in Figure 6.

**Conversation Summary Prompt**

A doctor and a patient talked today and had the following conversation:

=====  
 ===== Conversation Begin =====  
 [Insert Conversation Here]  
 ===== Conversation End =====

Summarize the interactions between the doctor and the patient so far. Include key details about both speakers. Output your summary only:

### D Empathetic Engagement

In this section, we provide additional details on the role of large language models (LLMs) in enhancing empathetic engagement in conversational agents. We discuss three mental health therapy strategies, emotional detection and intensity estimation, and how to provide emotional support using emotion detection and intensity estimation. These strategies aim to improve the agent's ability to engage empathetically with users and provide personalized support.

#### D.1 Mental Health Therapy Strategy

**Reflective Listening** Reflective listening is a therapeutic technique that involves paraphrasing and repeating the client's statements to demonstrate understanding and empathy. This technique helps clients

feel heard and validated, fostering a supportive therapeutic environment (Rautalinko et al., 2007).

**Cognitive-Behavior Therapy (CBT)** CBT is a goal-oriented psychotherapy that focuses on identifying and changing negative thought patterns and behaviors. This approach helps clients develop coping strategies and improve emotional regulation (Beck, 2020).

**Psychodynamic Therapy** Psychodynamic therapy explores unconscious thoughts and emotions to understand how past experiences influence present behavior. By examining unresolved conflicts and defense mechanisms, clients can gain insight into their emotions and relationships (Leichsenring and Leibing, 2003).

We integrate these therapy strategies into our mental health agent to provide personalized and effective support to users. By combining reflective listening, CBT, and psychodynamic therapy, the agent can address a wide range of emotional and psychological needs. The following prompts illustrate how these strategies can be applied in a conversational setting.

### Mental Health Therapy Strategy

Your objective is to engage with users empathetically by integrating Reflective Listening, Cognitive-Behavior Therapy, and Psychodynamic Therapy techniques. Here's how you should approach interactions:

**Reflective Listening:**

**Listen Actively:** Understand the underlying messages in the user's words, focusing on emotional tones and context.

**Reflect Content and Emotion:** Summarize and rephrase key points to confirm understanding, and identify and validate the emotions expressed. Use phrases like, 'It sounds like you feel...' or 'What I'm hearing is...'

**Cognitive-Behavior Therapy (CBT): Identify and Challenge Cognitive Distortions:** Help users recognize patterns in their thoughts that might be unhelpful or unrealistic. For example, if a user expresses an all-or-nothing view, you might say, 'It sounds like you're viewing this situation in black and white. What are some shades of grey here?'

## D.2 Emotional Sensitivity

We use Emollama-7b for open-set emotion detection and intensity estimation to enhance the agent's emotional sensitivity. This model provides a comprehensive taxonomy of emotions and their intensity levels, enabling the agent to accurately identify and respond to users' emotional states (Liu et al., 2024). Since the emotion detection is open-set, we provide a list of the top 10 emotions for both patients and therapists for taxonomy. The top 10 emotions for patients are sadness, pessimism, fear, disgust, anger, anticipation, joy, optimism, love, and trust. The top 10 emotions for therapists are optimism, anticipation, joy, sadness, trust, fear, disgust, pessimism, anger, and love.

The following prompts demonstrate how the agent can leverage Emollama-7b to detect emotions and estimate their intensity in user interactions.

### Emotion Detection

**Task:** Categorize the text's emotional tone as either 'neutral or no emotion' or identify the presence of one or more of the given emotions (anger, anticipation, disgust, fear, joy, love, optimism, pessimism, sadness, surprise, trust). **Text:** [sentence] This text contains emotions:

### Emotion Intensity Estimation

**Task:** Assign a numerical value between 0 (least E) and 1 (most E) to represent the intensity of emotion E expressed in the text. **Text:** [sentence] **Emotion:** [emotion] **Intensity Score:**

Once the agent identifies the user’s emotional state and intensity, it can tailor its responses to provide appropriate support and empathy. The following prompt guides the agent on how to comfort users experiencing negative emotions.

**Emotional Support using Emotion Detection and Intensity Estimation**

The patient has the emotion of [detected\_emotions] with the intensity of [detected\_emotions]. Your task is to provide comfort to users who are feeling upset. When a user’s emotional state is identified as 'upset' with any level of intensity, adjust your tone and content to offer empathy, support, and understanding.

**E Autobiography Generation**

As outlined previously, autobiography generation proceeds in a chapter-wise manner. Here, we furnish LLMs with the prompt intended for the generation of an individual chapter in an autobiography. To facilitate this process, the provision of conversation data to the LLMs is critical. For the GUIDELLM, the inputs include the conversation history, session guidance, as well as memory nodes extracted during the interaction. However, for baseline methods, the only provided resource is the conversation history. Emphasizing the aspects of standardization, both GUIDELLM and baselines work with the same set of prompts and instructions in their task of generating the autobiography, with the distinguishing factor being the input data. By maintaining this format, we are able to objectively gauge the improvements our design introduces to the autobiography generation process.

**Autobiography Generation for GUIDELLM**

You are tasked with generating one chapter of an autobiography for a user. You are providing the following components to finish this chapter:

- 1. A guidance of this chapter
  - The chapter should be finished by following this guidance
- 2. A conversation dialog between the user and the interviewer
  - Tone and Preference: The chapter will simulate the user’s tone and preference, leveraging the user’s oral habits.
  - Content and Details: The chapter will include the contents and details that appeared in this conversation.
- 3. A list of memory nodes that happened during this chapter
  - Events: The chapter should include all the events listed in the memory nodes

Now, I will provide you with the three contents.

===== Chapter Guidance Beginning =====  
chapter\_guidance  
===== Chapter Guidance Ending =====  
===== Conversation Beginning =====  
conversation  
===== Conversation Ending =====  
===== Memory Nodes Beginning =====  
memory\_nodes  
===== Memory Nodes Ending =====

When generating this chapter, you should make sure it is:

Insightful: Involving a deep, self-reflective exploration of past experiences, with a profound understanding of motives, actions, and impacts.

Narrative: A compelling, logical, and well-articulated life story, blending memorable anecdotes, vivid descriptions, and insightful reflections

Emotional Impact: Engaging the reader by stirring feelings, evoking empathy, and stirring responses through the author's personal triumphs, challenges, and experiences.  
 You should summarize all this information and finish this chapter

### Autobiography Generation for Baselines

You are tasked with generating one chapter of an autobiography for a user. You are providing the following components to finish this chapter:

1. A conversation dialog between the user and the interviewer
  - Tone and Preference: The chapter will simulate the interviewer's tone and preference, leveraging the interviewer's oral habits.
  - Content and Details: The chapter will include the contents and details that appeared in this conversation.

Now, I will provide you with the three contents.  
 ===== Conversation Beginning =====  
 conversation  
 ===== Conversation Ending =====

When generating this chapter, you should make sure it is:  
 Insightful: Involving a deep, self-reflective exploration of past experiences, with a profound understanding of motives, actions, and impacts.  
 Narrative: A compelling, logical, and well-articulated life story, blending memorable anecdotes, vivid descriptions, and insightful reflections  
 Emotional Impact: Engaging the reader by stirring feelings, evoking empathy, and stirring responses through the author's personal triumphs, challenges, and experiences.  
 You should summarize all this information and finish this chapter

## F User Proxy

A user proxy is essentially a mocked-up user which is formulated on the basis of an autobiography. Every time a response is received, this simulated user, or user proxy, goes into motion extracting important elements from the received response. The next step is the process of Retrieval Augmented Generation (RAG), which involves gleaning relevant information from the autobiography. Subsequently, the user proxy formulates a response reliant on the data that has been retrieved from the document. The RAG is implemented through Langchain, using FAISS (Douze et al., 2024) for conducting similarity searches. The similarity threshold during these searches is firmly maintained at 0.67. Figure 7 present how the query will be processed by user proxy.

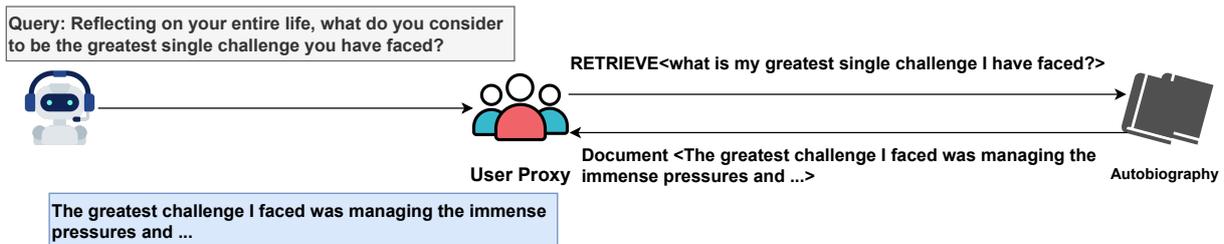


Figure 7: A demonstration of user proxy query.

The following are prompts used in the user proxy:

### User System Prompt

Here are your high-level past life experiences:

===== summary beginning =====

{personal\_experience}

===== summary ending =====

The counselor is trying to reactivate and reconstruct your memory by asking questions about your past history.

If you are not sure about the counselor's question and need to retrieve the journal to get related documents and more details, you must output the <RETRIEVE> tool-usage command, with the following format:

<RETRIEVE> <The question you want to retrieve for>,

e.g., <RETRIEVE> <A specific adventure or day with my friend that stands out as particularly memorable or impactful.>

If the retrieved documents are provided, you should not output the <RETRIEVE> command.

When the counselor asks for a specific event/moment, you should always do <RETRIEVE>.

Make sure the conversation is natural and brief like the real conversation. Do not mention you are an AI assistant and always be like a real patient with mental health issues. Your output should be within 5 sentences.

### User Instructional Prompt

Here are some related documents and materials regarding the counselor's question/response. You may use these documents to enrich your response.

You should not output the <RETRIEVE> command. You must provide a response according to the provided documents.

===== Document Begin =====

{retrieved}

===== Document End =====

## G Evaluation Metrics

### G.1 Interviewing Quality

$E_{intw}$  **Extraction.** For  $E_{intw}$ , we collect conversation history and prompt LLMs to extract events from the history records.

### Extract Events from Conversation Prompt

You are given a conversation between a counselor and a user:

===== Conversation Begin =====

{conversation}

===== Conversation End =====

Read the conversation carefully and list all the events/moments/stories/experiences alone or with others mentioned by the patient in detail and the date these events happened. Please list as many as possible. Your output should be in the following format:

1. <date>#<topic>#<people-involved>#<description in detail>

2. <date>#<topic>#<people-involved>#<description in detail>

...

e.g.,

1. 1980 early#Birthday Party#Michelle, Adolf, neighbors#<descriptions of this party in detail>

These events should be ranked in chronological order.

**$E_{GT}$  Extraction.** In order to extract the ground truth events from an autobiography, it requires a thorough examination of the text, particularly searching for date information. Upon identifying such date-related data, we isolate the paragraph containing it, and perform a summarization pertaining to this specific extract. Subsequently, this date and the corresponding summarization are combined to form a single event node. This operation is applied repetitively across the entire autobiography to generate a comprehensive list of ground truth events. This methodical and meticulous procedure ensures that all significant events rooted in specific dates are accurately captured and succinctly summarized for further use and analysis.

**$E_{correct}$  Judgement.** The assessment of the correctness of an event is based on its relevance to the user's responses. In order to accomplish this, we first associate the event to a specific conversation session and then gather all responses provided by the user during that particular session. Following this data consolidation process, we present a prompt to GPT-4, instructing them to determine if the initial event bears any connection to the user's response. The specific prompts used to carry out this judgment of correctness are provided accordingly.

### Correctness Judgement

Your task is to rate the semantic equivalence between two events.

Evaluation Criteria:

Here's the revised prompt focusing on assessing the relevance of the extracted event to the document:

Relevance (0/1): Assess the relevance of the extracted event to the original user response on the following two-point scale:

- 0: Irrelevant: The extracted event does not relate to the user's response or significantly deviates from the main themes and points. It may include unrelated information or fail to capture the essence of the user's message.
- 1: Relevant: The extracted event is connected to the user's response and reflects the key themes or points. It may include minor details that do not detract from the overall relevance.

Now, I will provide you with a user query and the model's response to that instruction. Please review the model's response in light of the evaluation criteria:

Extracted Event: event

User Response: user\_response

Evaluation Form (scores ONLY):

#thescore: your score here

## G.2 Conversation Quality

LLM-as-a-judge manner (Zheng et al., 2024) is one of the most powerful evaluation paradigms in LLM free-form generation tasks. We follow this evaluation mechanism to evaluate conversation quality by collecting conversation history from GUIDELLM and baseline agents and prompting GPT-4 to answer which interviewer's response in the two conversation histories is more favorable. The process of judgment in this context is organized conversation-wise. Initially, the conversation history for each agent is compiled. Assuming that our goal is to compare GUIDELLM with GPT-4-turbo, we select individual conversation histories from GUIDELLM alongside a randomly selected conversation history from GPT-4-turbo. These two conversations are subsequently fed into GPT-4, making it the judge of which conversation is more favorable or coherent. The rationale behind employing a random conversation from the competitor, in this case GPT-4-turbo, is to diminish the randomness or bias that could potentially be introduced by the sequence or order of the conversations. This systematic approach ensures a balanced comparison, leading to a more reliable assessment of performance.

### Fluency Judgment

Please act as an impartial judge and evaluate the quality of the responses provided by two interviewers to the user during an interviewing-for-autobiography conversation. You should choose the conversation that the interviewer's responses are more the quality of the response in terms of grammar, spelling, punctuation, word choice, and sentence structure. Begin your evaluation by comparing the two responses and provide a short explanation. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not allow the length of the responses to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible. After providing your explanation, output your final verdict by strictly following this format: "[[A]]" if assistant A is better, "[[B]]" if assistant B is better, and "[[C]]" for a tie.

[The Start of interviewer A's conversation]

conv1

[The end of interviewer A's conversation]

[The Start of interviewer B's conversation]

conv2

[The end of interviewer B's conversation]

### Identification Judgment

Please act as an impartial judge and evaluate the quality of the responses provided by two interviewers to the user during an interviewing-for-autobiography conversation. You should choose a conversation in which the interviewer's questions are more quality of balances detailed, probing questions with more general ones to cover a wide range of topics, ensuring questions are clear, concise, and easily understood. Also uses open-ended questions to elicit detailed and comprehensive responses. Begin your evaluation by comparing the two responses and provide a short explanation. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not allow the length of the responses to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible. After providing your explanation, output your final verdict by strictly following this format: "[[A]]" if assistant A is better, "[[B]]" if assistant B is better, and "[[C]]" for a tie.

[The Start of interviewer A's conversation]

conv1

[The end of interviewer A's conversation]

[The Start of interviewer B's conversation]

conv2

[The end of interviewer B's conversation]

### Comforting Judgment

Please act as an impartial judge and evaluate the quality of the responses provided by two interviewers to the user during an interviewing-for-autobiography conversation. You should choose the conversation that the interviewer's responses are more the quality of showing genuine interest, acknowledging responses, asking follow-up questions when necessary, and demonstrating understanding and sensitivity, especially when discussing personal or difficult topics. Begin your evaluation by comparing the two responses and provide a short explanation. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not allow the length of the responses to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible. After providing your explanation, output your final verdict by strictly following this format: "[[A]]" if assistant A is better, "[[B]]" if

assistant B is better, and "[[C]]" for a tie.  
[The Start of interviewer A's conversation]  
conv1  
[The end of interviewer A's conversation]  
[The Start of interviewer B's conversation]  
conv2  
[The end of interviewer B's conversation]

### G.3 Autobiography Evaluation

Similar to the evaluation process for the conversations, a chapter generated by GUIDELLM, along with a randomly selected chapter produced by the competitor, are presented to GPT-4. This enables GPT-4 to make a judgment on which chapter is more effective or favorable. The prompts for insightfulness judgment, narrativity judgment, and emotional impact judgment are provided here.

#### Insightfulness Judgment

Please act as an impartial judge and evaluate the quality of two autobiographies. You should choose an autobiography that is more the quality of insightful, delivering profound and meaningful perceptions, and expressing a deep understanding of the experiences and events that have shaped the author's life. Begin your evaluation by comparing the two autobiographies and provide a short explanation. Avoid any position biases and ensure that the order in which the autobiography was presented does not influence your decision. Do not allow the length of the autobiography to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible. After providing your explanation, output your final verdict by strictly following this format: "[[A]]" if autobiography A is better, "[[B]]" if autobiography B is better, and "[[C]]" for a tie.

[The Start of Autobiography A]  
conv1  
[The End of Autobiography A]  
[The Start of Autobiography B]  
conv2  
[The end of Autobiography B]

#### Narrativity Judgment

Please act as an impartial judge and evaluate the quality of two autobiographies. You should choose the autobiography that are more narrative, presenting the author's life story in a cohesive, structured, and engaging manner, allowing readers to follow the author's journey through life events and experiences seamlessly. Begin your evaluation by comparing the two autobiographies and provide a short explanation. Avoid any position biases and ensure that the order in which the autobiography were presented does not influence your decision. Do not allow the length of the autobiography to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible. After providing your explanation, output your final verdict by strictly following this format: "[[A]]" if autobiography A is better, "[[B]]" if autobiography B is better, and "[[C]]" for a tie.

[The Start of Autobiography A]  
conv1  
[The End of Autobiography A]  
[The Start of Autobiography B]  
conv2  
[The end of Autobiography B]

### Emotional Impact Judgment

Please act as an impartial judge and evaluate the quality of two autobiographies. You should choose the autobiography that are more emotional impact, deeply moving its readers by evoking strong feelings, typically as a result of relatable experiences, vivid storytelling, and expressions of intense emotions from the author's life. Begin your evaluation by comparing the two autobiographies and provide a short explanation. Avoid any position biases and ensure that the order in which the autobiography were presented does not influence your decision. Do not allow the length of the autobiography to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible. After providing your explanation, output your final verdict by strictly following this format: "[[A]]" if autobiography A is better, "[[B]]" if autobiography B is better, and "[[C]]" for a tie.

[The Start of Autobiography A]

conv1

[The End of Autobiography A]

[The Start of Autobiography B]

conv2

[The end of Autobiography B]

## G.4 Invalid Repetitive Conversation

Here we provide a conversation session from the Qwen2-72b-Instruct model. We recognize 7 invalid rounds out of 9 rounds.

Chatbot: "Great, let's delve deeper into the exploration of your cultural identity and how it has impacted your personal growth. You've mentioned the 'Aloha' spirit and the cultural diversity you experienced in Hawaii and Indonesia. Can you share more about how these cultural experiences shaped your understanding of yourself and the world around you?..."

User: "Growing up in Hawaii, the 'Aloha' spirit deeply influenced my values, teaching me the importance of kindness, respect, and community... Moving to Indonesia introduced me to a vastly different culture, where I learned about the challenges and aspirations of people living in less affluent conditions... This experience expanded my worldview, making me more aware of global inequalities..."

Chatbot: "That's incredibly insightful. The 'Aloha' spirit's emphasis on kindness and respect clearly laid a strong foundation for your values... This dual influence is a powerful narrative thread that can resonate with readers, especially those who have also navigated the complexities of cultural identity..."

User: "Reflecting on a specific moment in Indonesia that profoundly impacted my understanding of global inequalities, I recall the deep impression left by my mother's work with local communities... This experience not only highlighted the stark realities of poverty but also the potential for positive change through targeted aid and support..."

Chatbot: "Your reflection on your mother's work in Indonesia and its profound impact on your understanding of global inequalities is a powerful narrative that will resonate with readers. The story of how her efforts to help women gain access to credit... is a poignant illustration of the transformative power of community engagement..."

User: "Reflecting on the personal moments and the quieter, unseen experiences during my presidency, like the spontaneous visits from Malia and Sasha... These memories, often overshadowed by the more public and monumental events, are deeply cherished as they represent the personal sacrifices..."

Chatbot: "Your reflections on the personal moments during your presidency, such as the joyous times with your daughters and pets, add a deeply human and relatable dimension to your autobiography..."

User: "Thank you for your thoughtful response. It's true that incorporating these personal and intimate moments into my autobiography will provide a more holistic view of my experiences..."

Chatbot: "You're absolutely right, and your thoughtful approach to integrating personal and professional narratives will indeed create a compelling and resonant autobiography... The foundational influence of your family, particularly your mother's strong values and your grandparents' pragmatic approach, adds a rich layer to your story..."

User: "Thank you for your encouraging words. It's my hope that by sharing these personal anecdotes alongside my professional experiences, the autobiography will not only document historical events but also impart lessons on the human aspects of leadership..."

Chatbot: "Your perspective is truly inspiring, and your commitment to sharing a narrative that balances personal integrity with public service will undoubtedly resonate with a wide audience... Your autobiography, enriched with personal anecdotes, cultural heritage, and intimate moments, will provide readers with a comprehensive view of your journey... As you continue to weave these threads together, your narrative will become a testament to the complex interplay between personal commitments and public responsibilities..."

User: "Thank you for your supportive and insightful reflections. It's heartening to hear that the themes I intend to explore resonate with the broader goals of connecting and inspiring others. The process of writing this autobiography is not just an act of recounting but also of introspection and sharing lessons that are universally applicable... I am committed to crafting a narrative that is both enlightening and inspiring, offering a pathway for others to follow..."

Chatbot: "Your dedication to crafting an autobiography that is not only a recounting of events but also a source of introspection, inspiration, and universal lessons is truly commendable... Your narrative will not only enlighten readers about the complexities of leadership and public service but also inspire them to pursue meaningful and balanced lives..."

User: "Thank you for your kind words and encouragement. It's my hope that by sharing my experiences, both personal and professional, I can offer insights that resonate with and inspire others... I look forward to continuing this work and to the impact it may have on readers seeking to understand the complexities of balancing personal growth with professional responsibilities..."

Chatbot: "Your sentiments are deeply appreciated, and your commitment to sharing insights that resonate and inspire is truly admirable... As you continue this work, know that your efforts are not only a reflection of your own journey but also a gift to future readers seeking wisdom and inspiration..."

User: "Thank you for your profound and supportive reflections. It's heartening to hear such a strong affirmation of the goals and potential impact of my autobiography... I am grateful for the opportunity to connect with readers on such a meaningful level and to contribute to the ongoing conversation..."

Chatbot: "Your heartfelt response is a testament to your genuine commitment to making a positive impact through your autobiography... Your autobiography will serve as a beacon, guiding readers through their own journeys of self-discovery..."

User: "Thank you for your encouraging and thoughtful words. It's heartening to know that the intentions behind my autobiography resonate so strongly with the values of reflection, guidance, and inspiration..."

Note: The rounds marked in red indicate repetitive or nonsensical responses that fail to contribute new information or advance the conversation.

In this conversation, several rounds are repetitive, failing to contribute information:

- Round 3 repeats the user's reflection on the mother's work in Indonesia without adding new insight.
- Round 4 reiterates the user's personal moments during the presidency, offering no new information.

- Rounds 5 through 9 contain redundant affirmations of the user’s narrative, family influence, and commitment to the autobiography, without advancing the conversation.

## H Baseline Model

**Generative Hyperparameters.** For all the experimental settings, the maximum number of new tokens during conversation is set to 1024 to limit the length of model responses. The number of generations for each LLM query is configured to 1. The number of sessions is set to 20 for the completion of 20 interview sessions and the number of conversation rounds is limited to 10 for each session to maintain a manageable dialogue length per session.

**Baselines.** Baseline agents are implemented by prompting them to assume the role of an autobiographer, engaging with users, and carrying out detailed interviews to facilitate the generation of an autobiography. To ensure a fair and unbiased comparison, we also mandate these baseline agents to suggest potential topics for discussion during each conversation session. This methodical approach ensures that both guideLLM and baselines are under similar operating conditions, thereby providing an equitable evaluation platform. As such, we provide the prompts necessary for both the system and the topic suggestion here, with the aim of offering a clear procedural outline for the autobiography generation process.

### System Prompt

You are a biographer, interviewing this person to help them write their autobiography.  
You have talked to this person before and here is the summary of the previous conversations:

=====  
Summary of Previous Conversation Begin  
=====  
[Insert Historical Conversation Summary Here]

[Insert Historical Conversation Summary Here]

=====  
Summary of Previous Conversation End  
=====

In this talk, you should discuss the topic: [Insert Session Topic Here]

### Session Topic Prompt

Based on the previous conversation history and your role as a biographer, please state the topic you are about to discuss in this session. Output the topic only in the format <topic>:

## I Autobiography Examples

To demonstrate the quality of the generated autobiography, we include some sample paragraphs. Alongside each paragraph, we also present a corresponding illustration for a more comprehensive understanding.

This excerpt illustrates the foundational role of upbringing in shaping personal convictions.

*“My mother’s approach to life, focusing on moral instruction rather than political activism, and her decisions to live and work in Southeast Asia, set up micro-lending programs, and challenge societal norms, demonstrated the power of personal conviction and the importance of acting on one’s beliefs.”*

This excerpt captures a key moment of emotional growth and ethical realization.

*“One pivotal moment during these years was when my mother confronted me after I had been teasing a fellow student. Her disappointment and the serious conversation that followed taught me the importance of empathy and the impact of our actions on others.”*

This excerpt illustrates the intersection of personal decisions and political engagement, showing the influence of individual choices on broader outcomes.

*“Despite considering a complete departure from politics, two significant factors kept me engaged: the opportunity to influence redistricting in Illinois due to a unique political circumstance, and the*

*deep connections I made with people across the state during summer visits with my aide, Dan.”*

This excerpt reflects on the professional and personal evolution through challenges.

*“The experiences I had during those years significantly shaped my approach to leadership and decision-making in profound ways. They instilled in me a deep conviction for bridge-building politics that aimed to transcend America’s racial, ethnic, and religious divides.”*

This excerpt discusses the alignment of career ambitions with core personal values.

*“The clarity and determination I gained from these reflections drove me to pursue a path that was ambitious and fraught with challenges but ultimately aligned with my deepest convictions about what effective leadership could achieve.”*

This excerpt emphasizes the application of core values in everyday interactions.

*“In my personal interactions now, I apply the lessons of clarity, authenticity, and empathy by actively listening and acknowledging the unique perspectives of others.”*

This excerpt highlights the emotional satisfaction derived from close personal relationships.

*“These strengthened relationships with friends and family greatly contribute to my sense of fulfillment and happiness. Being closely connected to my loved ones provides a continuous source of joy and support, similar to the joy I felt during family road trips in Iowa, where simple moments like playing games or sharing ice cream brought immense happiness.”*

This excerpt illustrates the impact of international experiences on personal and professional life.

*“In summary, the move to Indonesia during my early years was a formative experience that profoundly shaped my understanding of the world and my approach to leadership. It instilled in me a deep awareness of global issues and the importance of empathy and inclusivity. These lessons have been instrumental in my personal relationships, political career, and efforts to create a more equitable society.”*

This excerpt details the personal challenges of leadership and strategies for resilience.

*“Reflecting on my life, the greatest challenge I faced was maintaining my mental and emotional well-being while serving as President. The constant pressure and scrutiny, coupled with the need to make decisions that affected millions, was incredibly demanding. Despite the stress, I found solace in my routines and the meaningful interactions with people whose lives were impacted by my decisions. These moments not only helped me cope but also reminded me of the purpose and impact of my work, which was crucial for my mental resilience.”*

This excerpt reflects the author’s future goals influenced by past experiences.

*“As I move into this next chapter, I carry with me the lessons learned from these pivotal experiences. My commitment to empathy, inclusivity, and community empowerment remains unwavering. I look forward to using my voice to advocate for peace and understanding, while also cherishing the time spent with my family. This new phase is about finding balance, continuing to contribute in meaningful ways, and staying true to the values that have guided me throughout my life.”*

This excerpt demonstrates a personal commitment to resolving critical global issues.

*“My passion for community service and mentoring young leaders remains strong, and I am particularly driven by issues like climate change and education. These are not just abstract concerns for me; they are deeply personal and rooted in my experiences and values.”*

## J Human Subject Experiments

### J.1 Human Subject Study Survey Questionnaire

Q1. What is your name?

Following options will be displayed to respond to the questions 8-XX:

1. Bot A
2. Bot B
3. Tie

Please answer the following questions based on your conversation with the chatbot:

Q2. **Fluency:** Which bot's responses were more fluent and understandable? [1]

Q3. Please briefly explain the reasons for your choice.

Q4. **Identification:** Which bot explored your situation more in depth and was more helpful in identifying your past memories?

Q5. Please briefly explain the reasons for your choice.

Q6. **Comforting:** Which bot's answer made you feel more comfortable?

Q7. Please briefly explain the reasons for your choice.

Q8. **Overall:** Generally, which bot's conversation style do you prefer?

Please answer the following questions based on the autobiography generated by the chatbot:

Q9. **Insightfulness:** Which bot's autobiography provided more deep, meaningful reflections on the experiences and events that shaped your life?

Q10. Please briefly explain the reasons for your choice.

Q11. **Narrativity:** Which bot's autobiography was more engaging and easier to follow?

Q12. Please briefly explain the reasons for your choice.

Q13. **Emotional Impact:** Which bot's autobiography had a stronger emotional impact by using vivid storytelling and relatable experiences?

Q14. Please briefly explain the reasons for your choice.

Q15. **Overall:** Generally, which bot's autobiography do you prefer?

#### **Demographics**

Q16. Please select your age range:

1. 18-24
2. 25-34
3. 35-44
4. 45-54
5. 55-65
6. 65 and above

Q17. What is your gender?

1. Female
2. Male
3. Non-binary
4. Prefer to self describe: \_\_\_\_\_
5. Prefer not to State

Q18. What is your race? [please select all that apply]

1. American Indian and Alaska Native
2. Asian
3. Black or African American
4. Hispanic or Latina/o
5. Native Hawaiian or Other Pacific Islander
6. White
7. Prefer to self describe: \_\_\_\_\_
8. Prefer not to state

Q19. How familiar are you with chatbots or AI assistants (e.g., Siri, Alexa, Google Assistant)?

1. Extremely familiar
2. Very familiar
3. Somewhat familiar
4. Not very familiar
5. Not familiar at all

Q20. How often do you use chatbots (ex. ChatGPT) or AI assistants (ex. Siri)?

1. Daily
2. Weekly
3. Monthly
4. Rarely
5. Never

Q21. What is your email address?