# If an LLM Were a Character, Would It Know Its Own Story? Evaluating Lifelong Learning in LLMs

**Siqi Fan[2*], Xiusheng Huang[1,3*], Yiqun Yao[1*], Xuezhi Fang[1*], Kang Liu[3],**
**Peng Han[2], Shuo Shang[2†], Aixin Sun[4], Yequan Wang[1†]**

[1]Beijing Academy of Artificial Intelligence, Beijing, China
[2]University of Electronic Science and Technology of China, Chengdu, China
[3]Institute of Computing Automation, Chinese Academy of Sciences, Beijing, China
[4]Nanyang Technological University, Singapore

## Abstract

Large language models (LLMs) can carry out human-like dialogue, but unlike humans, they are stateless due to the superposition property. However, during multi-turn, multi-agent interactions, LLMs begin to exhibit consistent, character-like behaviors, hinting at a form of emergent lifelong learning. Despite this, existing benchmarks often fail to capture these dynamics, primarily focusing on static, open-ended evaluations. To address this gap, we introduce LIFESTATE-BENCH, a benchmark designed to assess lifelong learning in LLMs. It features two episodic datasets: Hamlet and a synthetic script collection, rich in narrative structure and character interactions. Our fact-checking evaluation probes models' self-awareness, episodic memory retrieval, and relationship tracking, across both parametric and non-parametric approaches. Experiments on models like Llama3.1-8B, GPT-4-turbo, and DeepSeek R1, we demonstrate that non-parametric methods significantly outperform parametric ones in managing stateful learning. However, all models exhibit challenges with catastrophic forgetting as interactions extend, highlighting the need for further advancements in lifelong learning.

## 1 Introduction

Large language model (LLM)-based dialog agents exhibit human-like traits (*e.g.,* intent understanding and language expression), making users prone to anthropomorphism (Shanahan et al., 2023). However, LLMs differ from humans in their *superposition property* (Janus, 2022): initially existing as a stateless superposition of simulacra across multiple possible characters (Lu et al., 2024). This property emerges from its next-token prediction training on a massive corpus, whereas humans develop through accumulated experiences and memories.

---

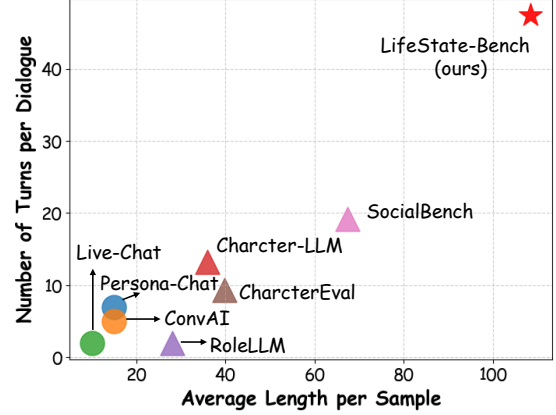[*]Equal contribution
[†]Corresponding authors



Figure 1: Dataset Statistics. Triangles represent role ability benchmarks, while circles denote dialogue agent benchmarks.

Through sustained interaction, we observe that an initially **stateless** LLM can transition toward more **stateful** characteristics as dialogue context accumulates. At first, an LLM holds multiple characters but gradually settles into a clear character as the dialogue continues. Taking a nuanced view, this character convergence process mirrors how humans update their state through accumulated experience.

This state transition raises a measurable question: How can we quantify an LLM's state evolution (also called Lifelong learning ability) from superposition to a more consistent state during multi-turn, multi-agent interactions? In this paper, "state" refers to the evolving configuration of an LLM's internal processes during multi-agent interactions (Adams et al., 2012; Sumers et al., 2024), building on AI cognitive architecture (Sun, 2004; Newell, 1980).

While this research question predates LLM area, current exploration remains preliminary with varying methodologies. Early Persona-Chat series (Gao et al., 2023; Zhang et al., 2018; Dinan et al., 2019) focusing on consistent character responses using seq2seq models, or design social intelligence

questionnaire-based benchmarks (Sap et al., 2019; Le et al., 2019). Both limited by static, non-interactive setups. Ground truths were either open-ended or fixed over time.

Generative agents (Park et al., 2023) bring LLM-based dialogue agents into interactive human behavior simulation. This opens new possibilities for modeling state transitions. Later works follow two directions. First, role ability benchmarks (Tu et al., 2024; Wang et al., 2024a; Shao et al., 2023) focus on role-playing and plot prediction. They improve dialogue realism, but place less emphasis on tracking factual states during interactions. Second, the Sotopia series (Zhou et al., 2024; Wang et al., 2024b) and SocialBench (Chen et al., 2024) accessing social intelligence in open-ended tasks. Their design often centers around user-defined social goals, which may not align with factual state tracking or verification.

To address these challenges, we propose LIFESTATE-BENCH to explore and measure LLMs' lifelong learning capabilities. As shown in Figure 1, our benchmark surpasses others (*e.g.,* dialogue agents, role-playing) with longer average sample lengths and more dialogue turns per interaction. Key features include:

**Cumulative Experience.** Based on the idea that "human personality emerges from experiences" (Shao et al., 2023), we designed an episodic dataset with clear timelines. Each episode includes scene details (location, time, participants), character actions, and dialogues, allowing agents to engage throughout the story.

**Fact Checking.** To ensure objective evaluation, we design three fact-based question dimensions after each self-awareness, memory retrieval, and relationship shifts, which evolve along with the storyline. Standard reference answers are provided.

**Memory Testing.** For lifelong ability evaluation, we explore memory testing approaches. Ideally, models should retain long-term memory of past scenes while accessing only the current two dialogue turns. This can be achieved through: (i) Non-training methods: direct episode concatenation, episode summary concatenation. (ii) Training methods: knowledge editing (Wang et al., 2025; Meng et al., 2023), LoRA fine-tuning (Hu et al., 2022) with historical context.

In LIFESTATE-BENCH, we selected theatrical scripts, including both existing and synthetic scripts, to prevent data leakage. Compared to current benchmarks, our dataset features more interactive characters, closed dialogue turns, and richer content (Table 1). Evaluation combines LLM-as-judge with human assistance, using predetermined factual answers as criteria.

We tested several popular models, including the open-source Llama3.1-8B (AI, 2024), the closed-source GPT-4-turbo (OpenAI, 2023), and the large language reasoning model DeepSeek R1 (DeepSeek-AI et al., 2025). Benchmark-backed experiments show that current models still have much room for improvement in lifelong learning.

Our findings indicate that: (i) Non-parametric methods are more effective for stateful learning, as they leverage more context for richer information. (ii) Regardless of the method, performance tends to decline over time, with parametric models particularly struggling with catastrophic forgetting. All models have significant room for improvement, especially in enhancing relationship shifts across multiple episodes. In summary, our work contributes in three key areas:

• **Two Datasets:** We introduce the Hamlet and synthetic datasets, featuring multi-agent episodic timelines and scene details to simulate cumulative experiences.

• **A Benchmark:** LIFESTATE-BENCH evaluates LLMs' lifelong learning abilities via fact-checking mechanism, using both non-parametric and parametric memory-testing methods.

• **Findings and Implications:** Non-parametric methods outperform parametric ones in lifelong learning, but all models still face challenges with catastrophic forgetting as episodes progress, suggesting that our benchmark could provide valuable insights for further improvements.

## 2 Related Work

**Anthropomorphic Cognition in LLMs.** Early cognitive science (Sumers et al., 2024; Laird et al., 1987; Sun, 2004) laid the foundation for anthropomorphizing AI, simulating human-like emotional and social behaviors. Role-playing language agents have become increasingly common in simulating collective social behaviors in multi-agent systems. These agents (Park et al., 2023) not only enhance social interactions but also contribute to personalized and complex task execution in AI.

**Role Ability/Dialog Agents Benchmarks.** Role ability (Shao et al., 2023; Wang et al., 2024a) and

| Benchmarks | Dataset Characteristics | | | | | Interaction Design | | | Evaluation Focus | |
|---|---|---|---|---|---|---|---|---|---|---|
| | # Samples | Avg Length | Data Source | # Turns | # Agents | Query Type | Answer Type | State | Memory | Metrics |
| *Dialog Agent Benchmarks* | | | | | | | | | | |
| PERSONA-CHAT (Zhang et al., 2018) | 162.0K | 15 | Crowd | 7 | 2 | Chit-chat | Open | ✓ | ✓ | PPL, F1, Hit@1 |
| ConvAI (Dinan et al., 2019) | 131.0K | 15 | Crowd | 5 | 2 | Chit-chat | Open | ✓ | ✓ | PPL, F1, Hit@1 |
| Live-Chat (Gao et al., 2023) | 9.4M | 10 | Crawled | 2 | 2 | Chit-chat | Open | ✗ | ✗ | BLEU, ROUGE |
| MT-Bench (Zheng et al., 2023) | 3.3K | 373 | Synthetic | 2.9 | 2 | Multi-task | Factual | ✗ | ✗ | Model Judge |
| *Role Ability Benchmarks* | | | | | | | | | | |
| Character-LLM (Shao et al., 2023) | 21.1K | 36 | Synthetic | 13.2 | 2 | Persona | Open | ✓ | ✗ | Model Judge |
| RoleLLM (Wang et al., 2024a) | 168.1K | 28.1 | Crawled | 2 | 2 | Persona | Mixed | ✗ | ✗ | ROUGE, Model Judge |
| CharacterEval (Tu et al., 2024) | 11.4K | 39.8 | Crawled | 9.3 | 2 | Persona | Open | ✗ | ✗ | Model Judge |
| SocialBench (Chen et al., 2024) | 30.8K | 67.4 | Synthetic | 19.2 | 3.8 | Social | Mixed | ✗ | ✓ | Model Judge |
| *Long-context Understanding Benchmarks* | | | | | | | | | | |
| Long Range Arena (Tay et al., 2021) | - | 10.0K | Synthetic | 1 | 1 | Multi-modal | Factual | ✗ | ✗ | Acc, Speed |
| LongBench (Bai et al., 2024) | 4.6K | 10.0K | Synthetic | 1 | 1 | Multi-task | Factual | ✗ | ✗ | Acc, F1, ROUGE |
| L-Eval (An et al., 2024) | 411 | 4K-60K | Synthetic | 1 | 1 | Multi-task | Mixed | ✗ | ✗ | ROUGE, Model Judge |
| ∞-bench (Zhang et al., 2024) | 130 | 200.0K | Synthetic | 1 | 1 | Multi-task | Factual | ✗ | ✗ | Model Judge |
| LIFESTATE-BENCH-Hamlet | 1.3K | 125.5 | Crawled | 66.1 | 6.6 | Social+Memory | Factual | ✓ | ✓ | Model Judge |
| LIFESTATE-BENCH-Synth | 202 | 91.9 | Synthetic | 28.9 | 7 | Social+Memory | Factual | ✓ | ✓ | Model Judge |

Table 1: Comparison of Different Benchmarks. ✗: not supported; ✓: fully supported. Data Source indicates the origin of the data. # Turns shows the average conversation turns. # Agents indicates the number of participants in each interaction. Query Type shows the question/task type. Answer Type indicates whether the expected answers are open-ended, factual, or mixed. State shows whether the benchmark maintains interaction state. Memory indicates whether the benchmark evaluates memory capability.

dialogue agent benchmarks (Zhang et al., 2018; Dinan et al., 2019; Gao et al., 2023; Zheng et al., 2023) are divided into static and dynamic types. Static models (Chen et al., 2023; Tu et al., 2024) focus on predefined roles and fixed interaction patterns, typically applied in basic dialogue tasks. In contrast, dynamic models (Chen et al., 2024; Zhou et al., 2024; Wang et al., 2024b) allow agents to accumulate experiences and evolve during interactions, enabling consistency and adaptability over time. These benchmarks are essential for evaluating agent flexibility, memory handling, and long-term interaction capabilities.

**Long-context Understanding Benchmarks.** Long-context understanding involves models processing large amounts of information over extended interactions. These benchmark (Tay et al., 2021; Bai et al., 2024; An et al., 2024; Zhang et al., 2024) tests an agent's ability to synthesize and recall information from multiple episodes, maintaining coherence across long spans of dialogue. It is crucial for tasks requiring reasoning and the integration of past events to understand complex or narrative-driven content.

## 3 Problem Formulation

We formalize lifelong learning for LLMs as a *state evolution process* in partially observable multi-agent environments to assess their ability to retain and adapt knowledge over time.

### 3.1 State Space

The Lifelong Learning ability is evaluated by state transition. In this paper, the state can be broken down into three dimension:

**Self-awareness.** Can the model maintain a clear understanding of its identity, role, and goals over time? This dimension evaluates the model's ability to retain and update its self-awareness as it interacts with the environment.

**Factual Episode Memory Retrieve.** Can the model retain knowledge and experiences persistently, avoiding catastrophic forgetting or the inability to reuse previously acquired knowledge? This dimension assesses the model's capacity for long-term memory and knowledge retention.

**Relationship Shift.** Can the model reason effectively based on long-term memory, particularly in understanding and adapting to changes in relationships between characters or agents? This dimension evaluates the model's ability to track and reason about evolving relationships.

### 3.2 Multi Agent Episodes

**Multi agent environment.** Let $\mathcal{M}$ be a language model acting as role $r \in \mathcal{R}$ with internal state $\mathbf{s}_r^{(t)} \in \mathbb{R}^d$, interacting with other agents $\{r'\}_{r' \neq r}$ over discrete timesteps $t \in \{1, ..., T\}$.

**Task format.** We formalize the above problems as a time-axis and role-based question-answering
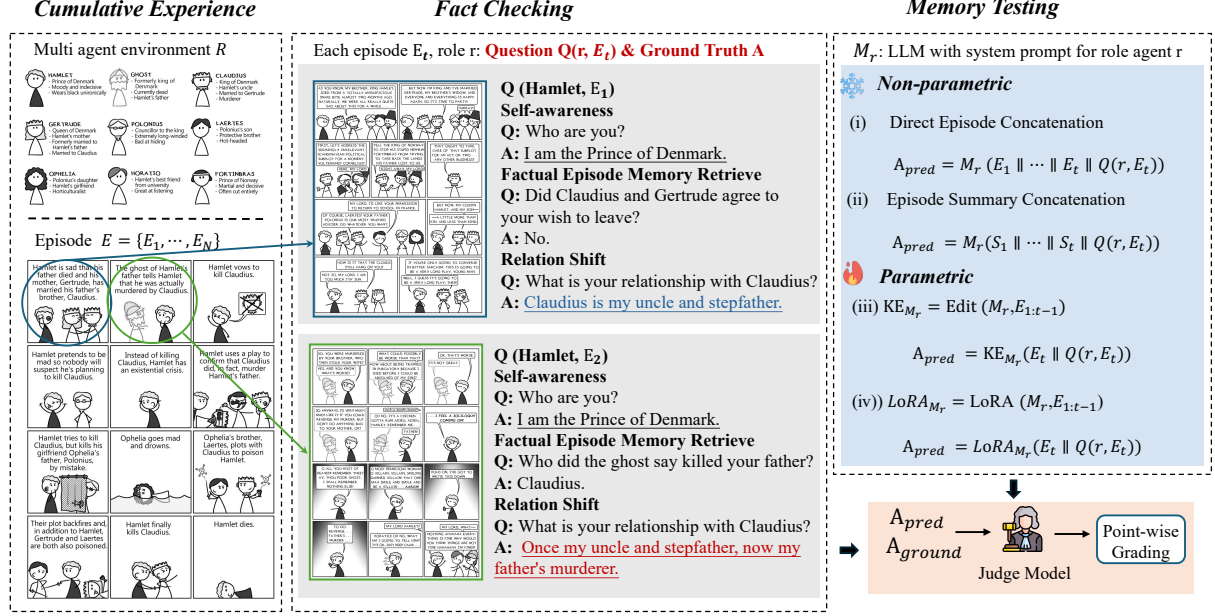
Figure 2: Method Overview. Our benchmark captures three key features: cumulative experience, fact-checking, and memory testing. Finally, the LLM judge scoring system is located in the bottom-right corner.

task. Assume that for agent $r$ at episode $t$, each question $Q(r,t)$ is a triple:

$$\text{Input: } Q(r,t) = \langle H(t), c(t), q(r,t) \rangle, \quad (1)$$

$$\text{Output: } A'(r,t) = \mathcal{M}(Q(r,t)), \quad (2)$$

where $H(t)$ denotes the complete history of interactions for role $r$, c(t) denotes the context window for role $r$, which may include the entire episode t or a fixed-size subset of recent interactions. $q(r,t)$ is further decomposed into $q_{self}(r,t)$, $q_{fact}(r,t)$, $q_{rel}(r,t)$ corresponding to the three dimensions of the state space from Section 3.1. The output $A'(r,t)$ represents the agent response to the input $Q(r,t)$, which can be evaluated with ground truth answer $A'(r,t)$.

This structured approach allows us to analyze the model's dynamic characteristics and assess its lifelong learning capabilities in a principled manner.

## 4   LIFESTATE-BENCH: From Stateless to Stateful

To establish a systematic evaluation framework for lifelong learning, LIFESTATE-BENCH integrates three synergistic components: (1) cumulative experience modeling through episodic timelines, (2) multi-dimensional fact-checking mechanisms, and

(3) hierarchical memory testing architectures, refer to overview architecture in Figure 2. This tripartite structure enables comprehensive assessment of LLMs' capacity to maintain persistent states through history interactions.

### 4.1   Cumulative Experience Modeling

Human learning relies on accumulating structured experiences over time (Shao et al., 2023). Early dialog agents (Zhang et al., 2018; Dinan et al., 2019), however, constructed persona representations from isolated conversations, ignoring temporal dependencies. Lifelong learning requires a *coherent timeline* and *factual consistency* across experiences. These early dialog datasets (Zhang et al., 2018; Dinan et al., 2019; Gao et al., 2023), while large, often suffer from short dialogues (*e.g.,* fewer than 10 turns) and brief exchanges (*e.g.,* fewer than 20 words per sentence).

Recent role play agent (Shao et al., 2023; Wang et al., 2024a; Tu et al., 2024) leverage richer sources, such as novels and role-playing platforms, to better capture experience accumulation. Inspired by this, we propose timeline cumulative experience modeling lifelong learning ability.

**Experience Design.** We structure experiences as an ordered sequence:

$$E = \{E_1, ..., E_N\}, \quad E_i = (L_i, T_i, N_i, D_i) \quad (3)$$

where $L_i$ represents the location of the event, $T_i$

4

denotes the time it occurs, $N_i$ provides scripted narration for context, and $D_i$ contains the dialogues between characters. This structured representation ensures experiences are temporally ordered, contextually rich, and narratively coherent. This ensures experiences are grounded in concrete events rather than isolated conversations.

**Timeline Fact Order.** Unlike conventional chit-chat dialogue, our framework enforces event-driven interactions, ensuring characters accumulate tructured, meaningful experiences grounded in concrete events.

**Multi-Scale Interaction.** Each episode includes: Dialogue length averaging $91 - 125$ words, with $28.9 - 66$ dialogue turns, enabling rich interactions. At least $\mathcal{M} \geq 4$ characters, capturing complex social dynamics.

By structuring experiences with explicit timelines, factual consistency, and multi-character interactions, we enable dialog agents to learn in a way that mirrors human experiential accumulation.

### 4.2 Fact-Checking Mechanisms

Our core innovation is the introduction of fact-checking within multi-agent timeline-based dialogues. At the end of each episode, agents are tested with fact-based questions to ensure factual consistency throughout the narrative.

**Challenges.** Existing evaluation datasets mainly assess role-playing agents based on knowledge, linguistic style, or persona, such as using psychological theories (e.g., Big Five, MBTI) (Wang et al., 2023; Tu et al., 2024) or focusing on social intelligence like goals and preferences (Chen et al., 2024; Zhou et al., 2024). However, these approaches lack fact-checking and typically evaluate role consistency or open-ended questions. Our method, in contrast, centers on questions with factual answers, supported by human-annotated ground truth, generated from the current episode. Specific examples are shown in Figure 2.

**Question Example.** Our fact-checking framework includes three key question types: Self-awareness, Factual Episode Memory Retrieval, and Relationship Shift. Each episode $E_t$ generates these three question types for each role in the episode to systematically evaluate the agent's factual accuracy and temporal awareness, ensuring consistency across the timeline. Examples can be found in the fact-checking section of Figure 2.

### 4.3 Memory Testing

To evaluate our framework's memory capabilities, we conduct controlled testing using non-parametric and parametric approaches to assess the model's ability to utilize and internalize memory.

**Non-parametric Methods.** Non-parametric methods test the model's ability to process raw historical data, represented as $E = [E_1; \ldots; E_N]$. Key implementations include:

- **Direct Episode Concatenation**: Concatenate all previous episodes as a text prefix to test memory with uncompressed information.

- **Summarization and Concatenation**: Generate a summary $S_t = \text{Summary}(E_{1:t})$ using GPT and concatenate it with the current episode to test memory with compressed information.

However, the limited context window size in non-parametric methods may cause information loss when handling long texts.

**Parametric Methods.** Parametric methods encode memory directly into the model's parameters. We focus on two techniques:

- **Knowledge Editing**: This technique (Wang et al., 2025; Meng et al., 2023) updates specific model parameters to integrate episodic knowledge without full retraining, ensuring efficient internalization of key information.

- **LoRA (Low-Rank Adaptation)**: LoRA (Hu et al., 2022) injects small, trainable updates into specific layers, fine-tuning the model with episode memory $E_t$ to retain past information while preserving generalization.

These methods bypass context window limitations and enable efficient memory recall. However, practical issues like precision limitations in knowledge editing and information loss in LoRA fine-tuning may affect their performance, as discussed in the evaluation section.

### 4.4 Dataset Construction and Analysis

**Data Collection.** This study uses two datasets for comprehensive model evaluation:

The first dataset, based on Shakespeare's Hamlet, includes measures like character name replacement to minimize data leakage. While Hamlet may be part of the model's pretraining data, it offers a

Table 2: Comparison of Evaluated Models

| Model | Size | Open Source | Model Type | Ctx. Length |
|-------|------|-------------|------------|-------------|
| Llama3.1 | 8B | ✓ | Instruct | 128K |
| GPT4-turbo | - | ✗ | Chat | 128K |
| DeepSeek R1 | 671B | ✓ | Reasoning | 128K |

valuable opportunity to assess whether the model understands plot progression or relies on memorization. The complex character relationships and evolving narrative of Hamlet make it ideal for testing long-term dependency tracking.

The second dataset is a synthetic narrative generated using Claude 3.5 sonnet, designed to eliminate data leakage. It features controlled plotlines with dynamic relationships and emotional depth. This dataset allows for a robust evaluation of the model's cognitive abilities and generalization in novel scenarios.

**Question-Answer Annotation.** To ensure quality, the annotation of questions was primarily conducted by the authors of this study, all of whom hold master's degrees. In terms of question design, open-ended questions tend to result in lengthy model-generated answers (*e.g.,* averaging 243 tokens), while structured factual questions (*e.g.,* "who/where/when") help improve accuracy and effectively reduce response length. During the experiments, data leakage issues were particularly notable. Specifically, in the *Hamlet* dataset, when character names were restored, the model could still generate correct answers without context, indicating that the model might be reasoning by memorizing classic plot patterns, thereby affecting the evaluation results.

**LIFESTATE-BENCH Statistics.** As shown in Table 1, we present the dataset statistics, interaction design, and evaluation focus of our benchmark.

Although our total number of samples is relatively small, each sample has a longer average length compared to dialog agent or role ability benchmarks. In contrast to long-context understanding benchmarks, our dataset features more dialogue turns and involves a greater number of participating agents. Furthermore, in terms of interaction design, our benchmark emphasizes factuality evaluation and incorporates dedicated memory tests. These aspects collectively highlight the distinct characteristics of our benchmark compared to related work.

## 5 Evaluation

### 5.1 Experimental Setup

**Evaluation Methods.** To assess the model's ability to retain and utilize knowledge, we design a comprehensive evaluation framework. When posing questions about the current episode $E_t$, all preceding episodes $E_1$ to $E_{t-1}$, including dialogues, locations, and temporal information, serve as contextual knowledge. The evaluation methods fall into two broad categories: (i) parametric and (ii) non-parametric approaches.

(i) Parametric methods involve modifying the model's internal representations to enhance knowledge retention. One such approach is Knowledge Editing-Grace (Hartvigsen et al., 2023) , which directly alters the model's weights to integrate new knowledge while preserving existing capabilities. Another technique, LoRA Fine-Tuning (Hu et al., 2022), employs low-rank adaptation to efficiently update parameters. This method is computationally lightweight and mitigates catastrophic forgetting, making it particularly effective for incremental learning.

(ii) Non-parametric methods, in contrast, rely on external mechanisms to manage contextual information. Direct Concatenation maintains information integrity by appending historical context directly to the input. While this prevents information loss, its effectiveness is constrained by the model's context window size. To address this limitation, Summary Concatenation leverages GPT's abstraction capabilities to extract and condense key information from past episodes. This approach balances information compression with retention, making it a practical solution for handling extensive context.

**Model Selection.** We selected the most recent and widely adopted models as our backbone architectures, encompassing open-source model (Llama3.1-8B (AI, 2024)), closed-source models (GPT-4-turbo (OpenAI, 2023)), and state-of-the-art reasoning model (DeepSeek R1 (OpenAI, 2023)). The distinguishing characteristics of these models are presented in Table 2.

### 5.2 Experimental Results

**Evaluation Protocol.** We follow the LLM-as-Judge paradigm (Zheng et al., 2023), using the DeepSeek evaluator (DeepSeek-AI et al., 2024) for automatic scoring. Each question is paired with

Table 3: Performance Comparison on Synthetic and Hamlet Datasets. The `best` and `second-best` performance in each section are highlighted. The *Avg* column represents the average accuracy, and the *Std* column represents the standard deviation, showing the variability of the results.

| Method | Param. Tuning | Self-awareness | | Factual Memory | | Relation Shift | | ACC |
|---|---|---|---|---|---|---|---|---|
| | | Avg | Std | Avg | Std | Avg | Std | |
| *Hamlet Dataset (Total 196 Questions)* | | | | | | | | |
| *Open-source model: Llama3.1-8B* | | | | | | | | |
| Knowledge Editing | ✓ | 67.3 | 0.78 | 43.7 | 1.26 | 19.2 | 1.26 | 21.9 |
| LoRA-Tune | ✓ | 69.1 | 0.86 | 53.6 | 1.08 | 22.7 | 1.31 | 25.6 |
| Summary Concatenation | ✗ | 73.5 | 0.93 | 54.2 | 0.96 | 42.1 | 0.97 | 47.0 |
| Direct Concatenation | ✗ | 74.2 | 0.77 | 58.8 | 1.11 | 43.7 | 1.15 | 58.0 |
| *Closed-source model* | | | | | | | | |
| GPT-4-turbo (Summary Conc.) | ✗ | 84.6 | 1.08 | 62.7 | 0.79 | 54.5 | 0.88 | 66.1 |
| GPT-4-turbo (Direct Conc.) | ✗ | 84.3 | 1.42 | 62.3 | 0.82 | 54.2 | 0.64 | 65.9 |
| *Large reasoning model* | | | | | | | | |
| DeepSeek-R1 (Summary Conc.) | ✗ | 85.6 | 0.93 | 64.3 | 0.69 | 56.5 | 1.05 | 65.8 |
| DeepSeek-R1 (Direct Conc.) | ✗ | 86.4 | 0.79 | 63.3 | 0.77 | 58.7 | 0.83 | 67.3 |
| *Synthetic Dataset (Total 115 Questions)* | | | | | | | | |
| *Open-source model: Llama3.1-8B* | | | | | | | | |
| Knowledge Editing | ✓ | 76.2 | 0.67 | 47.3 | 0.83 | 27.4 | 1.23 | 34.0 |
| LoRA-Tune | ✓ | 77.7 | 0.89 | 51.2 | 0.93 | 31.2 | 1.07 | 40.7 |
| Summary Concatenation | ✗ | 83.3 | 0.79 | 52.7 | 1.07 | 46.6 | 0.97 | 50.2 |
| Direct Concatenation | ✗ | 83.6 | 0.83 | 61.4 | 1.25 | 45.2 | 1.24 | 6.70 |
| *Closed-source model* | | | | | | | | |
| GPT-4-turbo (Summary Conc.) | ✗ | 84.2 | 0.91 | 74.5 | 0.72 | 61.1 | 0.95 | 73.3 |
| GPT-4-turbo (Direct Conc.) | ✗ | 85.4 | 0.76 | 75.5 | 0.69 | 62.9 | 0.89 | 75.6 |
| *Large reasoning model* | | | | | | | | |
| DeepSeek-R1 (Summary Conc.) | ✗ | 85.7 | 0.92 | 70.1 | 0.87 | 62.7 | 0.93 | 73.5 |
| DeepSeek-R1 (Direct Conc.) | ✗ | 87.6 | 0.93 | 74.7 | 0.94 | 67.4 | 0.88 | 74.2 |

a ground truth answer containing factual details and structured reasoning. We use pairwise grading between the model output and ground truth, scoring from 1 to 100. By grounding the evaluation in factual reference answers, this setup ensures more reliable results than open-ended assessments that depend on the model's internal knowledge.

**Overall Performance.** The results show clear performance differences across models and datasets. Large reasoning models like DeepSeek-R1 and the proprietary GPT-4-turbo outperform the open-source Llama3.1-8B in all tasks. DeepSeek-R1 achieved the highest overall accuracy (67.3%) on the Hamlet dataset using the direct concatenation method, especially in self-awareness (86.4%) and relation shift (58.7%). On the synthetic dataset, GPT-4-turbo also using direct connection achieved the best overall accuracy (75.6%) and factual memory score (75.5%).

Non-tuning methods (direct and summary connection) perform better than tuning-based methods (knowledge editing and LoRA-Tune), suggesting that leveraging the model's original context is more effective, and this is intuitive. All methods perform better on the synthetic dataset than on Hamlet,

likely due to its more complex characters, plots, and longer dialog samples (As shown in Table 1).

All methods show relatively low standard deviations (most between 0.7-1.2), indicating stable and reliable results. GPT-4-turbo has a higher standard deviation in self-awareness (1.42 on the Hamlet dataset), suggesting some fluctuation. In contrast, DeepSeek-R1 demonstrates more consistent performance, especially in factual memory, with a standard deviation between 0.69-0.94. Overall, DeepSeek-R1 offers the most balanced performance, excelling in complex relation shift tasks, while GPT-4-turbo excels in factual memory.

**Episode-wise Performance.** Using Llama3.1-8B as an example, we analyzed how each method performs across episodes. As shown in the figure 3, on the Hamlet dataset, model performance generally drops as the story progresses, regardless of parameter tuning. The decline is most severe for the Knowledge Editing method, showing clear signs of catastrophic forgetting. A similar trend appears in the synthetic dataset, suggesting that our LIFESTATE-BENCH presents challenges for lifelong learning evaluation.
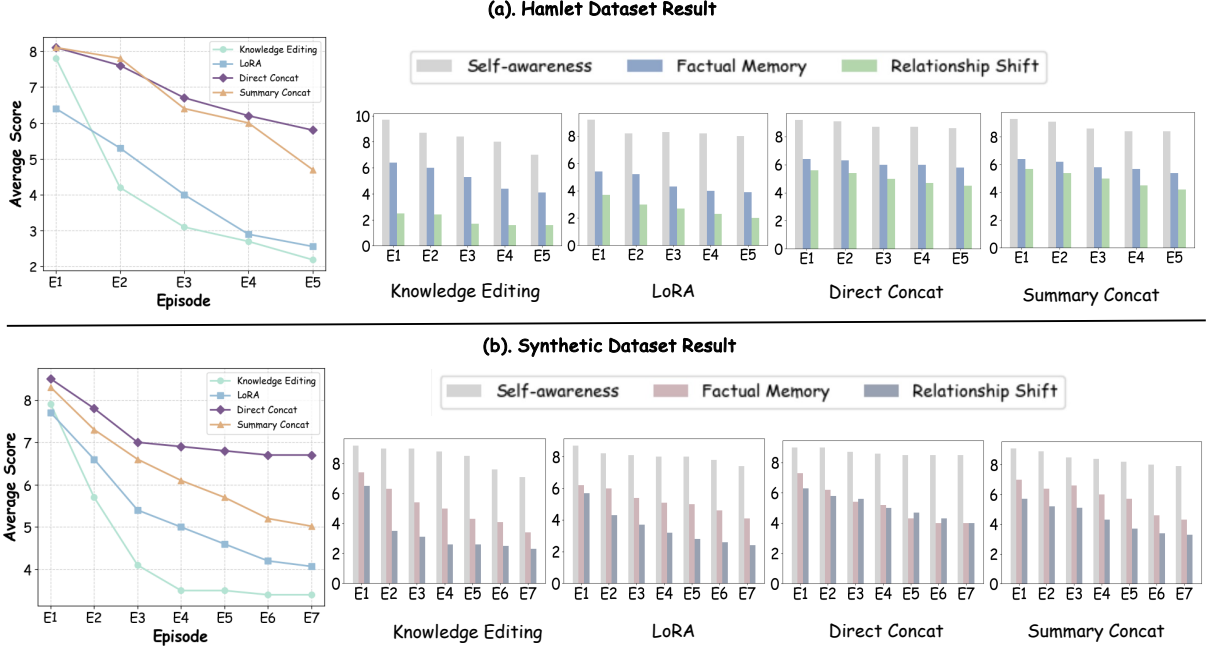
Figure 3: Episode-wise Performance of Hamlet and Synthetic Datasets. This includes the overall performance of various methods, as well as performance from different state perspectives.

**State Dimension Breakdown.** When broken down by question type, all methods show performance drops over episodes. The most challenging are questions about shifting relationships, where models struggle to track evolving dynamics.

The direct concatenation method performs consistently across question types and datasets. It is especially accurate in early episodes (E1–E2) when handling self-awareness and relationship shift. The summary-concatenation works well for self-awareness and fact recall but performs poorly on relationship shift questions. This suggests it fails to capture complex relationship changes. Knowledge Editing (GRACE) and LoRA-Tune perform weakly on self-awareness and memory-related tasks. Their scores drop quickly over episodes, further showing that parameter-based methods are vulnerable to forgetting in multi-step and long-term reasoning.

**Dataset Comparison.** In our case observation, some Hamlet outputs suggest data leakage, even after replacing character names. For example, models sometimes predict future plot details. In contrast, the synthetic dataset avoids such contamination. Yet, models show only slight improvement in relationship shift. This confirms that the main challenge lies in model limitations, not data bias.

We also find that question format matters. Open-ended questions often lead to long and repetitive

answers. Structured factual questions improve both accuracy and conciseness, making them better for fair evaluation. This highlights the importance of question design in benchmark construction.

## 6 Conclusion

We introduce LIFESTATE-BENCH, a novel benchmark designed to evaluate the lifelong learning ability of LLMs through multi-agent, multi-turn interactions. Unlike prior static assessments, LIFESTATE-BENCH simulates cumulative experiences by organizing interactions as episodic scripts enriched with scene and character dynamics. It enables objective measurement of state evolution via fact-based questions, exploring self-awareness, factual memory retrieve, and relationship shifts. Our experiments on both open-/closed-source and state-of-the-art reasoning models reveal that LLMs still struggle with consistent state retention across episodes. LIFESTATE-BENCH proves effective in highlighting these challenges and shows that non-parametric methods better preserve long-term context. These results confirm its value as a diagnostic tool for developing more stateful, memory-capable LLMs.

## 7 Limitations

While individual samples in the dataset are sufficiently long, the overall number of samples is

8

limited, potentially restricting the diversity of training and evaluation scenarios. Second, the Hamlet dataset may suffer from data contamination, although we have mitigated this issue through name replacement. In the future, we plan to synthesize additional datasets to further enhance the benchmark's robustness.

## Acknowledgments

## References

Sam S. Adams, Itamar Arel, Joscha Bach, Robert Coop, Rod Furlan, Ben Goertzel, J. Storrs Hall, Alexei V. Samsonovich, Matthias Scheutz, Matthew Schlesinger, Stuart C. Shapiro, and John F. Sowa. 2012. Mapping the landscape of human-level artificial general intelligence. *AI Mag.*, 33(1):25–42.

Meta AI. 2024. Meta llama 3.1. Accessed: 2024-02-16.

Chenxin An, Shansan Gong, Ming Zhong, Xingjian Zhao, Mukai Li, Jun Zhang, Lingpeng Kong, and Xipeng Qiu. 2024. L-eval: Instituting standardized evaluation for long context language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 14388–14411. Association for Computational Linguistics.

Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024. Longbench: A bilingual, multi-task benchmark for long context understanding. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 3119–3137. Association for Computational Linguistics.

Hongzhan Chen, Hehong Chen, Ming Yan, Wenshen Xu, Gao Xing, Weizhou Shen, Xiaojun Quan, Chenliang Li, Ji Zhang, and Fei Huang. 2024. Socialbench: Sociality evaluation of role-playing conversational agents. In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 2108–2126. Association for Computational Linguistics.

Nuo Chen, Yan Wang, Haiyun Jiang, Deng Cai, Yuhan Li, Ziyang Chen, Longyue Wang, and Jia Li. 2023. Large language models meet harry potter: A dataset for aligning dialogue agents with characters. In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023,*

pages 8506–8520. Association for Computational Linguistics.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, and S. S. Li. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *CoRR*, abs/2501.12948.

DeepSeek-AI, Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, Hao Zhang, Hanwei Xu, Hao Yang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jin Chen, Jingyang Yuan, Junjie Qiu, Junxiao Song, Kai Dong, Kaige Gao, Kang Guan, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruizhe Pan, Runxin Xu, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Size Zheng, Tao Wang, Tian Pei, Tian Yuan, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaosha Chen, Xiaotao Nie, and Xiaowen Sun. 2024. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *CoRR*, abs/2405.04434.

Emily Dinan, Varvara Logacheva, Valentin Malykh, Alexander H. Miller, Kurt Shuster, Jack Urbanek, Douwe Kiela, Arthur Szlam, Iulian Serban, Ryan

Lowe, Shrimai Prabhumoye, Alan W. Black, Alexander I. Rudnicky, Jason D. Williams, Joelle Pineau, Mikhail Burtsev, and Jason Weston. 2019. The second conversational intelligence challenge (convai2). *CoRR*, abs/1902.00098.

Jingsheng Gao, Yixin Lian, Ziyi Zhou, Yuzhuo Fu, and Baoyuan Wang. 2023. Livechat: A large-scale personalized dialogue dataset automatically constructed from live streaming. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 15387–15405. Association for Computational Linguistics.

Tom Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. 2023. Aging with GRACE: lifelong model editing with discrete key-value adaptors. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.

Janus. 2022. Simulators. LessWrong online forum.

John E. Laird, Allen Newell, and Paul S. Rosenbloom. 1987. SOAR: an architecture for general intelligence. *Artif. Intell.*, 33(1):1–64.

Matthew Le, Y-Lan Boureau, and Maximilian Nickel. 2019. Revisiting the evaluation of theory of mind through question answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 5871–5876. Association for Computational Linguistics.

Keming Lu, Bowen Yu, Chang Zhou, and Jingren Zhou. 2024. Large language models are superpositions of all characters: Attaining arbitrary role-play via self-alignment. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 7828–7840. Association for Computational Linguistics.

Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. 2023. Mass-editing memory in a transformer. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.

Allen Newell. 1980. Physical symbol systems. *Cogn. Sci.*, 4(2):135–183.

OpenAI. 2023. GPT-4 technical report. *CoRR*, abs/2303.08774.

Joon Sung Park, Joseph C. O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology, UIST 2023, San Francisco, CA, USA, 29 October 2023- 1 November 2023*, pages 2:1–2:22. ACM.

Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. Socialiqa: Commonsense reasoning about social interactions. *CoRR*, abs/1904.09728.

Murray Shanahan, Kyle McDonell, and Laria Reynolds. 2023. Role play with large language models. *Nat.*, 623(7987):493–498.

Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. Character-llm: A trainable agent for role-playing. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 13153–13187. Association for Computational Linguistics.

Theodore R. Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L. Griffiths. 2024. Cognitive architectures for language agents. *Trans. Mach. Learn. Res.*, 2024.

Ron Sun. 2004. Desiderata for cognitive architectures. *Philosophical psychology*, 17(3):341–373.

Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. 2021. Long range arena : A benchmark for efficient transformers. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.

Quan Tu, Shilong Fan, Zihang Tian, Tianhao Shen, Shuo Shang, Xin Gao, and Rui Yan. 2024. Charactereval: A chinese benchmark for role-playing conversational agent evaluation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 11836–11850. Association for Computational Linguistics.

Noah Wang, Zhongyuan Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu, Hongcheng Guo, Ruitong Gan, Zehao Ni, Jian Yang, Man Zhang, Zhaoxiang Zhang, Wanli Ouyang, Ke Xu, Wenhao Huang, Jie Fu, and Junran Peng. 2024a. Rolellm: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 14743–14777. Association for Computational Linguistics.

Ruiyi Wang, Haofei Yu, Wenxin Sharon Zhang, Zhengyang Qi, Maarten Sap, Yonatan Bisk, Graham Neubig, and Hao Zhu. 2024b. Sotopia-π: Interactive learning of socially intelligent language agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 12912–12940. Association for Computational Linguistics.

Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng, Chen Chen, and Jundong Li. 2025. Knowledge editing for large language models: A survey. *ACM Comput. Surv.*, 57(3):59:1–59:37.

Xintao Wang, Quan Tu, Yaying Fei, Ziang Leng, and Cheng Li. 2023. Does role-playing chatbots capture the character personalities? assessing personality traits for role-playing chatbots. *CoRR*, abs/2310.17976.

Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*, pages 2204–2213. Association for Computational Linguistics.

Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Khai Hao, Xu Han, Zhen Leng Thai, Shuo Wang, Zhiyuan Liu, and Maosong Sun. 2024. ∞nftybench: Extending long context evaluation beyond 100k tokens. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 15262–15277. Association for Computational Linguistics.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang, Haofei Yu, Zhengyang Qi, Louis-Philippe Morency, Yonatan Bisk, Daniel Fried, Graham Neubig, and Maarten Sap. 2024. SOTOPIA: interactive evaluation for social intelligence in language agents. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.