

# From *What to Respond* to *When to Respond*: Timely Response Generation for Open-domain Dialog Agents

Anonymous ACL submission

## Abstract

While research on dialog response generation has primarily focused on generating coherent responses conditioning on textual context, the critical question of *when to respond* grounded on the temporal context remains underexplored. To bridge this gap, we propose a novel task named timely dialog response generation and introduce the TIMELYCHAT benchmark to evaluate two key aspects: response timing and time-conditioned responses, which focus on the capabilities of language models to predict appropriate delays and delayed responses. Additionally, we construct a large-scale training dataset by leveraging unlabeled event knowledge from a temporal commonsense knowledge graph and employing a large language model (LLM) to synthesize 55K event-driven dialogs. We then train TIMER, a dialog agent designed to proactively predict time intervals and generate timely responses that align with those intervals. Experimental results show that TIMER outperforms instruction-tuned LLMs and other time-aware baselines in both turn-level and dialog-level evaluations. We publicly release our data, model, and code.<sup>1</sup>

## 1 Introduction

Open-domain dialog systems have introduced various aspects to improve human-likeness of their own, such as persona (Zhang et al., 2018; Ahn et al., 2023), long-term memory (Xu et al., 2022a,b), commonsense (Zhou et al., 2021; Qin et al., 2021), emotional support (Rashkin et al., 2019; Liu et al., 2021; Zhang et al., 2024), and role-play (Shao et al., 2023; Li et al., 2023). These efforts have led to the success of commercial chat services like Replika and Character AI, which have met the public’s demand for social companion chatbots (Chaturvedi et al., 2023; Guingrich and Graziano, 2023). The

<sup>1</sup><https://anonymous.4open.science/r/timelychat-EADE/>



Figure 1: An illustrative example of a timely dialog agent. Unlike delay-agnostic agents that can only provide instant responses, a timely dialog agent proactively predicts response delays as well as responses by considering the temporal context of the conversation, enabling human-like interactions.

pursuit for human-like chatbots still remains important with the remarkable advancements in large language models (LLMs), along with the growing demand for AI agents capable of engaging in more natural and human-like ways.

Research on dialog response generation has predominantly focused on generating appropriate and consistent next utterances, conditioning on the textual information within dialog contexts. Meanwhile, although the question of *what to respond* has received considerable attention, the problem of *when to respond* remains underexplored, despite its crucial role in enabling real-time dialog agents to appropriately ground their responses on the temporal contexts regarding the status of ongoing conversational events. For instance, as illus-

trated in Figure 1, if an agent generates only instant responses without considering response timing, it can cause repetitive interactions without conversational progress or produce responses that do not align with the temporal context of the conversational event. In contrast, by incorporating response timing, an agent can maintain a natural flow while providing timely responses. *This requires grounding responses on the temporal context tied to the status of the event, mirroring the way humans naturally adapt their responses in human-to-human conversations.*

However, it is inherently challenging to simulate such scenarios with dialog models trained on existing datasets. Most dialog datasets lack explicit temporal context and are created under the tacit assumption that interactions occur instantly. Additionally, collecting real-time conversations where temporal context is naturally embedded (e.g., text messages between individuals) is highly restricted due to privacy concerns and ethical considerations.

In this work, we propose a novel task named **Timely Dialog Response Generation**, which aims to generate not only coherent responses but also to consider the temporal context associated with ongoing events. Specifically, it focuses on predicting the necessary time interval for the next utterance and generating a corresponding time-conditioned response. We introduce TIMELYCHAT dataset and propose a benchmark to assess two key aspects: response timing prediction and time-conditioned response generation. To create diverse event-driven dialogs, we combine the human-annotated event-duration pairs from a temporal commonsense knowledge graph with the powerful dialog generation capability of an LLM.

Furthermore, we introduce a large-scale dataset comprising 55K event-driven dialogs for supervised fine-tuning (SFT). To address the challenges of costly and labor-intensive manual annotation, we utilize unlabeled event sources from a large-scale temporal commonsense knowledge graph and leverage an LLM to pseudo-label event durations and synthesize diverse event-driven dialogs. Using this dataset, we present TIMER, a dialog model fine-tuned with a multi-task learning objective that jointly predicts the time interval and generates the corresponding response.

Evaluation results on the proposed benchmark demonstrate that TIMER outperforms both instruction-tuned LLMs and dialog models trained on time-aware datasets in generating time-

conditioned responses and predicting time intervals consistent with temporal commonsense. Furthermore, in dialog-level evaluations, TIMER distinguishes between situations requiring delayed responses and those requiring instant responses more effectively, and generates more timely responses that align well with the predicted time intervals. We publicly release the TIMELYCHAT benchmark, training data, and our timely dialog agent named TIMER to facilitate further research in this area.

## 2 Related Work

Long-term dialog involves conversations that unfold over multiple sessions with time intervals between sessions. Xu et al. (2022a) introduce Multi-Session Chat (MSC), which consists of up to five sessions separated by certain time intervals, resembling interactions in messaging platforms. Jang et al. (2023) emphasize the significance of speaker relationships in long-term dialogs and propose Conversation Chronicles (CC), a large-scale LLM-generated dataset that incorporates a wider range of time intervals and fine-grained speaker information. Maharana et al. (2024) present Lo-CoMo, a very long-term dialog dataset covering up to 32 sessions, along with a benchmark designed to assess various long-term memory capabilities. However, prior research primarily focuses on recalling persona sentences or past events from previous sessions, without addressing the temporal context between ongoing events and time intervals in real-time conversations. A notable attempt to incorporate such relations is GapChat (Zhang et al., 2023), which introduces an event timeline to capture event progression over given time intervals. Our work moves beyond the assumption of predetermined time intervals and instead necessitates a proactive dialog agent capable of dynamically determining realistic time delays based on temporal context.

## 3 Task Definition

We introduce a new task named **Timely Dialog Response Generation**, which aims to generate contextually appropriate responses while incorporating temporal considerations from the dialog history. A key temporal factor that influences a response is how much time has passed since the previous utterance. To capture this, we define *time interval* as our primary temporal context, which represents the relative time difference (e.g., 10 minutes) between utterances. Formally, we model the conditional

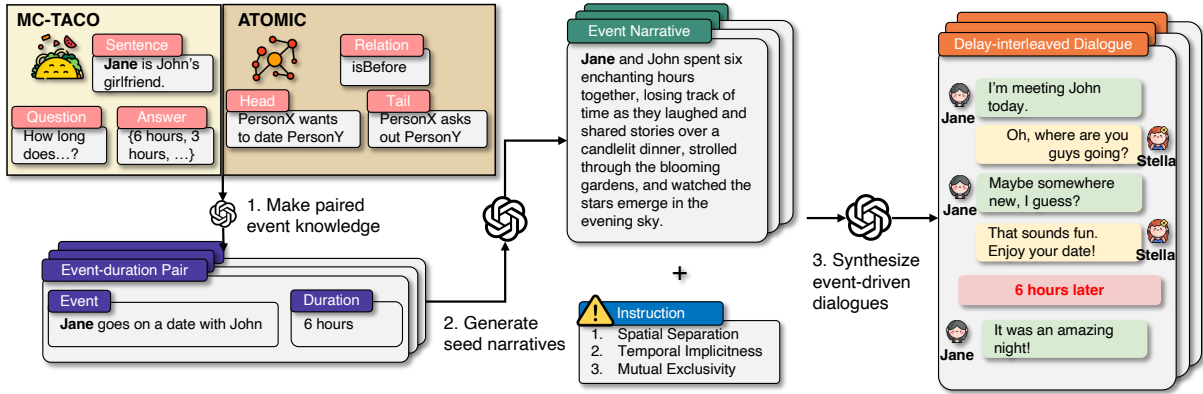


Figure 2: Overall process of data construction method. Two different knowledge sources represent the same example for better understanding. Note that due to constraints imposed by the instruction, Jane’s conversation partner becomes Stella, not John.

probability distribution  $P_\theta$  of a response  $r_t$  at  $t$ -th turn given the textual context  $u$  and the temporal context  $\tau$ :

$$r_t \sim P_\theta(u_t | u_1, \dots, u_{t-1}, \tau_2, \dots, \tau_t), \quad (1)$$

where  $u_t$  denotes the utterance at  $t$ -th turn and  $\tau_t$  denotes the elapsed time between  $u_{t-1}$  and  $u_t$ . This probability distribution can be further decomposed into two subtasks, which are the main focus of this work.

**Subtask 1. Response Timing Prediction** The first task is to predict the optimal timing for delivering messages to users. Mathematically, this involves predicting the  $t$ -th time interval  $\hat{\tau}_t$  given the available contexts:

$$\hat{\tau}_t \sim P_\theta(\tau_t | u_1, \dots, u_{t-1}, \tau_2, \dots, \tau_{t-1}). \quad (2)$$

**Subtask 2. Time-conditioned Response Generation** The subsequent task is to generate a contextually appropriate response while incorporating the predicted timing for message delivery:

$$r_t \sim P_\theta(u_t | u_1, \dots, u_{t-1}, \tau_2, \dots, \tau_{t-1}, \hat{\tau}_t). \quad (3)$$

Note that this task formulation challenges the widely held assumption that dialog agents should always respond to user messages instantly. Instead, it takes temporal context into account, i.e., the amount of elapsed time, to determine when a response should be generated.

## 4 TIMELYCHAT Benchmark

We construct TIMELYCHAT benchmark to assess the timely response generation capabilities of dialog models. To this end, we first craft high-quality

timely conversations through temporal knowledge base and LLMs and then design two evaluation processes. Figure 2 shows the overall construction process of our benchmark.

### 4.1 Data Construction

We incorporate temporal information into dialogs using a temporal commonsense knowledge base. This knowledge base captures various event-related temporal dynamics which is well suited for transforming temporal context into event-driven dialogs. By identifying temporal patterns, we seamlessly integrate them into conversations, utilizing the sophisticated dialog generation capabilities of LLMs.

**Event Knowledge Extraction.** We first obtain a rich and reliable source of daily events and their typical durations for crafting event-driven conversations with temporal context. To this end, we utilize the event duration category of MC-TACO dataset (Zhou et al., 2019). The dataset consists of sentences for specific events, queries to ask the typical duration of the event (e.g., "How long does it take to ...?"), and human-annotated ground-truth answers. We utilize the sentences with ground-truth answers, i.e., event-duration pair, to synthesize event-driven conversations. During data construction, we excluded the examples whose temporal intervals shorter than one minute or longer than 24 hours to simulate realistic temporal delay in daily dialog situations. Lastly, we instruct GPT-4o (Achiam et al., 2023) with these sentences and event-duration pairs to generate descriptive sentences. It integrates the event and its duration into coherent sentences, forming seed narratives for dialog generation.

**Timely Dialog Generation.** With the extracted temporal event knowledge, we instruct GPT-4o to generate conversations. Our instruction contains the conditions that the generated dialogs must satisfy:

- **Spatial Separation:** The scenario must involve one speaker experiencing an event while conversing with another speaker about it. This ensures there are no contradictions arising from both speakers being in the same spatial context.
- **Temporal Implicitness:** The response must avoid direct references to the elapsed time. This condition reduces the occurrence of dull responses that simply acknowledge the time interval and, more importantly, prevents lexical overlap with the ground-truth time interval, which could create shortcuts in the generation process.
- **Mutual Exclusivity:** The time-conditioned response must become untimely under contrary temporal conditions. In other words, a delayed response should be incoherent under an instant condition with no time interval, and an instant response should be incoherent when a time interval exists. It prevents generating time-agnostic responses that remain coherent regardless of the temporal context.

Along with these instructions, we provide one randomly selected example from six author-written dialogs, each ranging from 5 to 10 turns, to prevent ill-formed outputs and diversify dialog lengths. After manual inspection and filtering out low-quality dialogs that did not meet all the conditions, the final synthesized dataset consists of 324 dialogs, with an average length of 6.5 turns. All prompts and examples used in the construction process are provided in Appendix A.

## 4.2 Evaluation Protocols

With the crafted conversations, we propose two evaluation approaches to assess the abilities of dialog agents to generate timely responses: turn-level and dialog-level.

**Turn-level Evaluation.** In turn-level evaluation, we assess each subtask on the target response. For response timing prediction, a model predicts the time interval required for the next utterance given a dialog context. We then evaluate (1) whether the

model correctly classifies the next turn as either delayed or instant, and (2) how close is the predicted interval to the ground truth. We measure precision, recall, false positive rate (FPR), and F1 for the binary classification, and root mean squared logarithmic error (RMSLE) for regression by converting each time interval into minutes. For response generation, a model generates a time-conditioned response given a dialog context and ground-truth time interval. We measure BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and BERTScore (Zhang et al., 2020) as reference-based metrics. Additionally, we measure naturalness (Mehri et al., 2022) and time-specificity (Tsunomori et al., 2023) on a 5-point scale, adopting LLM-as-a-judge (Liu et al., 2023) for automatic evaluation.

**Dialog-level Evaluation.** One crucial quality of a timely dialog agent is its ability to introduce appropriate delays considering the temporal context while maintaining a natural conversational flow. Inspired by dialog-level evaluation methods with model-to-model interactions (Li et al., 2019; Zhou et al., 2024), we provide an event-driven scenario and let an agent converse with GPT-4o as a user simulator (Yoon et al., 2024; Kazi et al., 2024; Wang et al., 2024) for the fixed number of turns to measure dialog-level metrics. We measure coherence (Mehri et al., 2022) and dialog-level time-specificity to assess the quality of the agent’s responses, and measure delay appropriateness that considers both the timing and duration of delays, using LLM-as-a-judge with a 5-point scale. The evaluation criteria of LLM-as-a-judge metrics and simulator instructions are detailed in Appendix B.

## 5 TIMER: A Dialog Agent for Timely Responses

### 5.1 Training Data Augmentation with Unlabeled Knowledge

Utilizing paired event-duration knowledge is essential for creating conversations that simulate timely responses. However, manually constructing such annotations is both costly and labor-intensive, posing a challenge to creating large-scale datasets for training LMs. To overcome this limitation, we leverage unlabeled event knowledge graphs and harness the capabilities of GPT-3.5 to construct large-scale paired knowledge and generate synthetic dialogs. This approach significantly reduces the manual effort required while enabling the creation of extensive training data.

Dataset	# Sessions	Construction Method	Time Granularity	Event-grounded	# Events
MSC (Xu et al., 2022a)	13K (4.4K)	Crowdsourcing	hours - weeks	✗	-
CC (Jang et al., 2023)	1M (160K)	LLM-generated	hours - years	✗	-
LoCoMo (Maharana et al., 2024)	842 (-)	LLM-gen + Crowd	days - months	✗	-
GapChat (Zhang et al., 2023)	2.6K (782)	LLM-gen + Crowd	minutes - years	✓	128
<b>TIMELYCHAT (Ours)</b>	55K	LLM-generated	minutes - hours	✓	55K

Table 1: Comparison of long-term dialog datasets interleaved with time intervals. The number in parentheses under the # Sessions column represents the count of sessions with time intervals within a day. Event-grounded indicates whether the dialogs reflect the temporal context associated with events or not.

**Event Knowledge Extraction.** We extract event knowledge from the ATOMIC<sub>20</sub> dataset (Hwang et al., 2021), a large-scale commonsense knowledge graph containing the event-centered category represented as event triplets (i.e., head, relation, and tail), which capture diverse temporal dynamics. To make more natural dialogs, we randomly replace the anonymized person names (e.g., PersonX) in the triplets with common names of US SSN applicants, following the method by Kim et al. (2023). Subsequently, we prompt GPT-3.5 to integrate these triplets into single-sentence event descriptions, producing more natural and coherent event representations.

**Event Duration Estimation.** Since the event triplets in ATOMIC<sub>20</sub> do not include annotated durations, we utilize GPT-3.5 to estimate typical durations. Specifically, we provide GPT-3.5 with the event descriptions and prompt it to extract the main event and predict its typical duration, which is then used as a pseudo label. We filter out examples where the predicted duration is less than a minute or exceeds 24 hours.

**Dialog Generation with Bootstrap Examples.** We prompt GPT-3.5 using the instructions detailed in Section 4.1. During initial iterations, we observed that providing only the instructions often led to ill-formed dialogs, such as speaker mismatches or non-alternating turns. To address these issues and improve dialog quality, we include a one-shot demonstration sampled from the TIMELYCHAT set in each prompt. All prompts used in the construction process are presented in Appendix A.1.

The resulting dataset consists of 55K events paired with their corresponding dialogs. Compared to existing long-term dialog datasets in Table 1, our dataset includes a significantly larger amount of even-grounded dialogs without requiring costly human annotation and handles time intervals with finer granularity.

## 5.2 Time-augmented Training with Multi-task Learning Objectives

The goal of our training approach is to predict an appropriate time interval for delaying the response based on the temporal context of the conversation and then generate a time-conditioned response corresponding to the interval. For this purpose, we introduce a time interval prediction step before generating each turn’s utterance.

We propose a training approach for timely dialog response generation, as formalized in Eqs. 2 and 3. For each turn consisting of a speaker identifier and a text utterance, we insert a time interval. We prepend prefix tokens to distinguish each component, formatting the input as <SPK>  $s_t$  <TIME>  $\tau_t$  <UTT>  $u_t$ , where  $s_t$ ,  $\tau_t$ , and  $u_t$  denote the speaker, the time interval, and the utterance at the  $t$ -th turn, respectively. For turns within the dialog context, we set  $\tau = 0$ , indicating no delay, maintain coherence and align with typical instant responses.

From these inputs, we define two losses for training: response timing prediction loss and response generation loss. The losses are defined as follows:

$$\begin{aligned} \mathcal{L}_{\text{time}} &= -\frac{1}{N} \sum_{i=1}^N \sum_{t=2}^T \log p(\tau_t \mid s_{\leq t}, \tau_{<t}, u_{<t}), \\ \mathcal{L}_{\text{response}} &= -\frac{1}{N} \sum_{i=1}^N \sum_{t=2}^T \log p(u_t \mid s_{\leq t}, \tau_{\leq t}, u_{<t}), \end{aligned} \quad (4)$$

where  $N$  is the number of training examples, and  $T$  is the number of turns in a dialog.

The final multi-task learning objective is given as follows:

$$\mathcal{L} = \mathcal{L}_{\text{response}} + \lambda \mathcal{L}_{\text{time}}. \quad (5)$$

This approach ensures that the model learns both to predict appropriate time intervals and to generate time-conditioned responses effectively.

## 6 Experiments

### 6.1 Baselines

We evaluate two types of dialog agents for simulating timely dialog response generation: instruction-tuned models and time-aware models. The instruction-tuned models include LLMs optimized for dialog use cases. We select 8B and 70B models of LLaMA 3.1 Instruct (Dubey et al., 2024) as open-source chat models, and GPT-3.5 and GPT-4o as proprietary models. We experiment with zero-shot, few-shot, and chain-of-thought (CoT) (Wei et al., 2022) prompting strategies to investigate the effectiveness of in-context learning without task-specific fine-tuning. The time-aware models are trained on dialog datasets where time intervals are interleaved. We compare the following models:

- **MSC 3B** (Xu et al., 2022a): Trained on BlenderBot (Roller et al., 2021) using the MSC dataset, which includes time intervals between sessions.
- **ReBot 400M** (Jang et al., 2023): Trained on BART-Large (Lewis et al., 2020) using the CC dataset, which consists of large-scale LLM-generated dialogs.
- **GapChat 3B** (Zhang et al., 2023): Trained on MSC using the GapChat dataset, which incorporates event progress based on time intervals.

To prevent a judge LLM from introducing potential biases such as self-preference bias (Wataoka et al., 2024; Panickssery et al., 2024), we utilize Claude Sonnet 4.5 as the evaluator across both turn-level and dialog-level evaluations. Implementation details of all models including TIMER 3B are described in Appendix B.1.

### 6.2 Turn-level Evaluation Results

**Response Timing Prediction.** Table 2 presents the results of response timing prediction on the TIMELYCHAT. Overall, instruction-tuned models exhibit significantly low precision and F1 scores, and high FPR. This suggests that these models tend to over-predict the need for a delay, potentially introducing unnecessary intervals that disrupt the conversational flow. Although few-shot and CoT strategies slightly improve F1 scores across all LLMs, they sometimes negatively impact FPR compared to zero-shot prompting. In contrast, TIMER

Model	Precision $\uparrow$	Recall $\uparrow$	F1 $\uparrow$	FPR $\downarrow$	RMSLE $\downarrow$
<b>LLaMA 3.1 8B</b>					
<i>Zero-shot</i>	0.172	0.691	0.276	0.606	2.853
<i>Few-shot</i>	0.164	<b>0.920</b>	0.279	0.855	2.807
<i>CoT</i>	0.160	0.824	0.268	0.790	2.854
<b>LLaMA 3.1 70B</b>					
<i>Zero-shot</i>	0.153	0.735	0.254	0.740	2.479
<i>Few-shot</i>	0.197	0.602	0.297	0.448	2.326
<i>CoT</i>	0.194	0.836	0.315	0.636	3.066
<b>GPT-3.5</b>					
<i>Zero-shot</i>	0.143	0.784	0.241	0.861	2.763
<i>Few-shot</i>	0.212	0.349	0.264	0.237	2.146
<i>CoT</i>	0.186	0.762	0.299	0.609	2.667
<b>GPT-4o</b>					
<i>Zero-shot</i>	0.266	0.259	0.263	0.131	1.956
<i>Few-shot</i>	0.227	0.423	0.295	0.099	2.252
<i>CoT</i>	0.202	0.852	0.326	0.615	2.938
<b>TIMER 3B (Ours)</b>	<b>0.783</b>	0.799	<b>0.791</b>	<b>0.041</b>	<b>1.189</b>

Table 2: Results of response timing prediction. For few-shot and CoT strategies, we provide balanced 2-shot demonstrations which consist of one delayed example and one instant example, along with the task description used in zero-shot prompting.

3B achieves the highest F1 score and the lowest FPR compared to instruction-tuned models. Even the best-performing GPT-4o still lags significantly behind the fine-tuned TIMER 3B model.

Likewise, when it comes to predicting the length of time intervals, in-context learning methods fail to enhance performance effectively. While few-shot prompting achieves a lower RMSLE than CoT across all LLMs, it does not consistently outperform zero-shot prompting, as demonstrated by GPT-4o’s results. These findings indicate that prompting with task descriptions and demonstrations alone is insufficient to reliably predict whether to pose a delay and how long it should last. In contrast, task-specific fine-tuning is essential for effectively learning these capabilities.

**Time-conditioned Response Generation.** Table 3 shows the time-conditioned response generation performance on the TIMELYCHAT. For instruction-tuned models, we observe that zero-shot performance tends to improve as model size increases across all metrics. Among all LLMs, few-shot prompting consistently outperforms zero-shot prompting, while CoT prompting performs the worst in terms of naturalness and time-specificity. This aligns with previous findings that LLMs struggle to generate helpful CoT rationales for dialog response generation (Chae et al., 2023).

Meanwhile, models trained on existing dialog datasets that include time intervals exhibit poor overall performance. Notably, these models

Model	B-2	R-L	BS	Nat.	Spec.
INSTRUCTION-TUNED MODELS					
<b>LLaMA 3.1 8B</b>					
<i>Zero-shot</i>	5.38	12.38	86.21	2.24	2.78
<i>Few-shot</i>	7.63	13.47	86.85	2.81	2.83
<i>CoT</i>	6.17	12.50	86.23	1.31	1.84
<b>LLaMA 3.1 70B</b>					
<i>Zero-shot</i>	6.84	12.71	85.90	2.86	2.35
<i>Few-shot</i>	8.35	14.83	87.16	3.92	2.89
<i>CoT</i>	9.01	15.01	87.12	1.07	3.04
<b>GPT-3.5</b>					
<i>Zero-shot</i>	9.97	17.13	87.54	4.26	2.86
<i>Few-shot</i>	11.23	17.81	87.77	4.25	2.76
<i>CoT</i>	8.86	15.14	86.79	2.38	2.15
<b>GPT-4o</b>					
<i>Zero-shot</i>	9.17	16.76	87.35	4.75	3.09
<i>Few-shot</i>	11.15	18.51	87.91	<b>4.78</b>	3.09
<i>CoT</i>	10.25	17.13	87.52	3.21	2.84
TIME-AWARE MODELS					
MSC 3B	3.26	8.94	85.18	1.36	1.02
ReBot 400M	3.85	9.32	85.72	3.26	1.02
GapChat 3B	3.59	8.61	85.22	1.87	1.05
TIMER 3B (Ours)	<b>16.08</b>	<b>22.26</b>	<b>88.74</b>	3.84	<b>3.36</b>

Table 3: Results of time-conditioned response generation on TIMELYCHAT. B-2, R-L, BS, Nat., and Spec. refer to BLEU-2, ROUGE-L, BERTScore, naturalness, and time-specificity, respectively.

469 achieve low time-specificity, indicating that they  
470 struggle to generate timely responses conditioned  
471 on given time intervals. This stems from the char-  
472 acteristics that time intervals in existing long-term  
473 dialog datasets are assigned arbitrarily rather than  
474 based on the temporal context of ongoing events,  
475 making it difficult for models to learn the condi-  
476 tional distribution of responses based on the given  
477 interval. For example, we find that these mod-  
478 els frequently generate generic greeting messages,  
479 failing to capture the temporal nuances of timely  
480 responses. In contrast, TIMER 3B, despite having  
481 a smaller model size, achieves comparable natural-  
482 ness to instruction-tuned LLMs and even surpasses  
483 LLaMA 3.1 models. More importantly, it achieves  
484 the highest time-specificity, demonstrating that our  
485 training approach enables response generation that  
486 aligns well with event-specific temporal contexts.

### 487 6.3 Dialog-level Evaluation Results

488 Beyond turn-level evaluation, we also conduct  
489 dialog-level evaluation to assess whether a dialog  
490 agent can introduce temporally contextual delays  
491 at appropriate timing without disrupting the con-  
492 versational flow. We let the four zero-shot LLMs

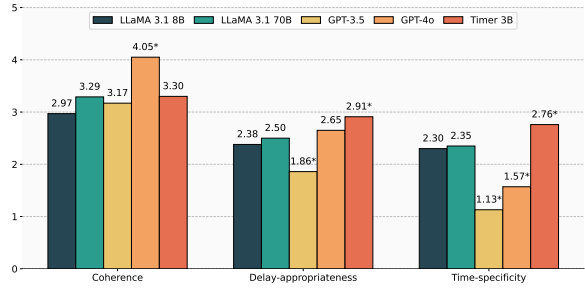


Figure 3: Results of simulated interactions for timely dialog agents. We perform pairwise t-test and denote statistically significant score differences from the other models with asterisk (\*) ( $p < 0.05$ ).

Setting	Metric	$\kappa$	Win	Tie	Loss
Turn-level	Naturalness	0.446	19%	24%	<b>57%</b>
	Time Specificity	0.508	<b>54%</b>	24%	22%
Dialog-level	Coherence	0.637	21%	18%	<b>61%</b>
	Delay Appropriateness	0.446	<b>40%</b>	23%	37%
	Time Specificity	0.496	<b>40%</b>	23%	37%

Table 4: Pairwise human evaluation results on both turn-level and dialog-level metrics. Win/Tie/Loss rates of TIMER 3B against zero-shot GPT-4o are presented. Fleiss'  $\kappa$  is also reported to measure inter-annotator agreements.

493 from the previous experiments, along with TIMER  
494 3B, engage in 10 interactions with the simulator  
495 described in Section 4.2. To simulate event-driven  
496 dialog, we provide the first turn of conversations  
497 from the TIMELYCHAT as the initial interaction.

498 We sample dialogs that include at least one de-  
499 layed response and report three dialog-level met-  
500 rics in Figure 3. GPT-4o achieves the highest co-  
501 herence among the models, demonstrating its abil-  
502 ity to maintain a natural conversation flow, while  
503 TIMER 3B achieves the second-highest coherence  
504 score. Notably, TIMER 3B shows the highest delay-  
505 appropriateness and time-specificity scores. This  
506 suggests that TIMER 3B effectively considers both  
507 dialog context and temporal context to predict de-  
508 lays with appropriate timing and duration. Addi-  
509 tionally, it generates delayed responses that are  
510 coherent only when a delay is given, thereby jus-  
511 tifying and necessitating the delay. In contrast,  
512 LLaMA 3.1 8B and 70B exhibit relatively lower  
513 delay-appropriateness, while GPT-3.5 and GPT-  
514 4o achieve lower time-specificity scores. We fur-  
515 ther analyze these findings in the case study in  
516 Appendix C.

Method	Timing Prediction		Response Generation		
	F1	RMSLE	R-L	Nat.	Spec.
<b>TIMER 3B</b>	<b>0.791</b>	<b>1.189</b>	<b>22.26</b>	<b>3.84</b>	<b>3.36</b>
w/o Special Tokens	0.006	1.699	17.09	2.75	2.76
Utterance First	0.249	3.451	12.07	1.57	1.81
w/o Time Loss	0.299	3.987	13.38	2.54	2.43

Table 5: Ablation study on each component of our training scheme. In Utterance First setting, we train a model to generate an utterance first and then predict the response timing conditioned on it.

## 6.4 Human Evaluation Results

To investigate the reliability of LLM-based evaluation, we also conduct human evaluations on both turn-level and dialog-level metrics. We recruit twenty graduate students as annotators, who are fluent in English but have no background on our work. We provide them with the same evaluation criteria used for LLM-based evaluation, and ask them to perform blinded comparisons between responses or dialog from two different models. Each example is evaluated three times by different annotators, and aggregated by majority voting. Table 4 presents the results for both turn-level and dialog-level metrics.

For turn-level metrics, we compare TIMER 3B with the most competitive baseline, zero-shot GPT-4o, on 200 randomly sampled examples. While TIMER 3B falls short of GPT-4o in terms of naturalness, it slightly outperforms GPT-4o in time-specificity, which is consistent with the LLM-based evaluation results observed in Table 3.

For dialog-level metrics, we compare TIMER 3B with zero-shot GPT-4o on 80 randomly sampled examples. Again, TIMER 3B lags behind GPT-4o in coherence, but it outperforms GPT-4o in both delay-appropriateness and time-specificity. This finding aligns with the results shown in Figure 3, indicating that the proposed evaluation criteria and LLM-based evaluation are reliable measures for assessing desired model behavior. We observe that the annotators show substantial ( $0.6 < \kappa \leq 0.8$ ) agreement on coherence, and moderate ( $0.4 < \kappa \leq 0.6$ ) agreements on the other metrics, demonstrating the reliability of human evaluation.

## 6.5 Ablation Study

To demonstrate the effectiveness of the proposed training scheme, we provide the impacts of each component on task performance in Table 5.

**Effects of Input Formats.** Special tokens (e.g.,  $\langle \text{UTT} \rangle$ ,  $\langle \text{TIME} \rangle$ ) play a crucial role for indicating which component of an utterance to generate. Although trained with both time intervals and responses, the model trained with no special tokens between the utterance components struggles to predict time intervals at the right place, and is not capable of predicting intervals other than 0 minutes. This results in arbitrary generation of delayed or instantaneous responses thereby causes poor controllability of the model.

The order of performing the two subtasks significantly affects the performance. When trained to generate a response first and then predict the interval conditioned on it, the model tend to predict delays even for instantaneous responses, resulting in poor performance on timing prediction. Notably, this format degrades response generation performance the most compared to the other settings, since the model is not provided with time conditions directly for generation.

**Effects of Timing Prediction Loss.** We investigate the effectiveness of our multi-task objective. When the timing prediction loss is removed ( $\lambda = 0$ ), the model struggles at learning the timing prediction itself, resulting in random prediction. Notably, response generation performance degrades compared to TIMER 3B, even though the model is still optimized with time-conditioned context inputs. This implies that the timing prediction objective helps improve response generation capability, demonstrating the necessity of the proposed training objective.

## 7 Conclusion

We highlighted the necessity for open-domain dialog agents to consider not only the response itself but also the timing of it based on the temporal context related to the conversational event. We formulated this challenge as the timely dialog response generation, and introduced the TIMELY-CHAT benchmark for turn-level and dialog-level evaluations. Additionally, we proposed a large-scale SFT dataset and a time-augmented training approach, which we used to train the TIMER 3B model that proactively predicts the time interval for the next utterance and then generate a time-conditioned response. TIMER 3B outperforms baseline models on the benchmark and demonstrates its ability to generate both appropriate time intervals and responses while maintaining natural flow.

## 8 Limitations

In this study, we predict event duration by mapping it to discrete values (e.g., 30 minutes). However, a more realistic assumption would be to consider it as a continuous time range (e.g., 2-6 hours). As future work, we aim to generalize this assumption to enable more fine-grained control over response delays. Additionally, while we use simulated dialogs with a few number of turns for dialog-level evaluation, further research could explore longer interactions across diverse social environments to analyze the correlation between human-likeness and user experience. Finally, beyond the fine-tuning and in-context learning methods used in our experiments, more task-specific training approaches could be developed to further enhance performance.

## 9 Ethics Statements

The proposed dataset was designed to assess capabilities related to response timing and time-conditioned response in event-driven conversations. To achieve this, we utilized event knowledge from publicly available datasets from various sources, and LLM-generated contents either without or with some modification if necessary. During this process, there is a possibility that harmful content or inappropriate biases existing in the original data may have been conveyed, or may have arisen due to limitations of filtering techniques. We reject any form of violence, discrimination, or offensive language, and our dataset and experimental results do not represent such values. If any harmful content or privacy infringement is identified within our data, we kindly request immediate notification to the authors. In the event of such cases being reported, we will apply the highest ethical standards and take appropriate actions.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. [Gpt-4 technical report](#). *arXiv preprint arXiv:2303.08774*.

Jaewoo Ahn, Yeda Song, Sangdoon Yun, and Gunhee Kim. 2023. [MPCHAT: Towards multimodal persona-grounded conversation](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3354–3377, Toronto, Canada. Association for Computational Linguistics.

Hyungjoo Chae, Yongho Song, Kai Ong, Taeyoon Kwon, Minjin Kim, Youngjae Yu, Dongha Lee, Dongyeop Kang, and Jinyoung Yeo. 2023. [Dialogue chain-of-thought distillation for commonsense-aware conversational agents](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5606–5632, Singapore. Association for Computational Linguistics.

Rijul Chaturvedi, Sanjeev Verma, Ronnie Das, and Yogesh K. Dwivedi. 2023. [Social companionship with artificial intelligence: Recent trends and future avenues](#). *Technological Forecasting and Social Change*, 193:122634.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. [The llama 3 herd of models](#). *arXiv preprint arXiv:2407.21783*.

Rose E Guingrich and Michael SA Graziano. 2023. [Chatbots as social companions: How people perceive consciousness, human likeness, and social health benefits in machines](#). *arXiv preprint arXiv:2311.10599*.

Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras, Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and Yejin Choi. 2021. [\(comet-\) atomic 2020: On symbolic and neural commonsense knowledge graphs](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(7):6384–6392.

Jihyoung Jang, Minseong Boo, and Hyoungun Kim. 2023. [Conversation chronicles: Towards diverse temporal and relational dynamics in multi-session conversations](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13584–13606, Singapore. Association for Computational Linguistics.

Taaha Kazi, Ruiliang Lyu, Sizhe Zhou, Dilek Hakkani-Tur, and Gokhan Tur. 2024. [Large language models as user-agents for evaluating task-oriented-dialogue systems](#). *Preprint*, arXiv:2411.09972.

Hyunwoo Kim, Jack Hessel, Liwei Jiang, Peter West, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Bras, Malihe Alikhani, Gunhee Kim, Maarten Sap, and Yejin Choi. 2023. [SODA: Million-scale dialogue distillation with social commonsense contextualization](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12930–12949, Singapore. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

711	Cheng Li, Ziang Leng, Chenxi Yan, Junyi Shen, Hao Wang, Weishi MI, Yaying Fei, Xiaoyang Feng, Song Yan, HaoSheng Wang, Linkang Zhan, Yaokai Jia, Pingyu Wu, and Haozhen Sun. 2023. <a href="#">Chatharuhi: Reviving anime character in reality via large language model</a> . <i>Preprint</i> , arXiv:2308.09597.	768
712		769
713		
714		
715		
716		
717	Margaret Li, Jason Weston, and Stephen Roller. 2019. <a href="#">Acute-eval: Improved dialogue evaluation with optimized questions and multi-turn comparisons</a> . <i>Preprint</i> , arXiv:1909.03087.	
718		
719		
720		
721	Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for automatic evaluation of summaries</a> . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	
722		
723		
724		
725	Siyang Liu, Chujie Zheng, Orianna Demasi, Sahand Sabour, Yu Li, Zhou Yu, Yong Jiang, and Minlie Huang. 2021. <a href="#">Towards emotional support dialog systems</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3469–3483, Online. Association for Computational Linguistics.	
726		
727		
728		
729		
730		
731		
732		
733		
734	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. <a href="#">G-eval: NLG evaluation using gpt-4 with better human alignment</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2511–2522, Singapore. Association for Computational Linguistics.	
735		
736		
737		
738		
739		
740		
741	Ilya Loshchilov and Frank Hutter. 2019. <a href="#">Decoupled weight decay regularization</a> . In <i>International Conference on Learning Representations</i> .	
742		
743		
744	Adyasha Maharana, Dong-Ho Lee, Sergey Tulyakov, Mohit Bansal, Francesco Barbieri, and Yuwei Fang. 2024. <a href="#">Evaluating very long-term conversational memory of LLM agents</a> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 13851–13870, Bangkok, Thailand. Association for Computational Linguistics.	
745		
746		
747		
748		
749		
750		
751		
752	Shikib Mehri, Jinho Choi, Luis Fernando D’Haro, Jan Deriu, Maxine Eskenazi, Milica Gasic, Kallirroi Georgila, Dilek Hakkani-Tur, Zekang Li, Verena Rieser, Samira Shaikh, David Traum, Yi-Ting Yeh, Zhou Yu, Yizhe Zhang, and Chen Zhang. 2022. <a href="#">Report from the nsf future directions workshop on automatic evaluation of dialog: Research directions and challenges</a> . <i>Preprint</i> , arXiv:2203.10012.	
753		
754		
755		
756		
757		
758		
759		
760	Paulius Micekevicius, Sharan Narang, Jonah Alben, Gregory Damos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, and Hao Wu. 2018. <a href="#">Mixed precision training</a> . In <i>International Conference on Learning Representations</i> .	
761		
762		
763		
764		
765		
766	Arjun Panickssery, Samuel Bowman, and Shi Feng. 2024. <a href="#">Llm evaluators recognize and favor their own generations</a> . <i>Advances in Neural Information Processing Systems</i> , 37:68772–68802.	776
767		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823

824	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. <a href="#">Chain-of-thought prompting elicits reasoning in large language models</a> . In <i>Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22</i> , Red Hook, NY, USA. Curran Associates Inc.	880
825		881
826		882
827		883
828		884
829		885
830		886
831	Jing Xu, Arthur Szlam, and Jason Weston. 2022a. <a href="#">Beyond goldfish memory: Long-term open-domain conversation</a> . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5180–5197, Dublin, Ireland. Association for Computational Linguistics.	887
832		888
833		
834		
835		
836		
837	Xinchao Xu, Zhibin Gou, Wenquan Wu, Zheng-Yu Niu, Hua Wu, Haifeng Wang, and Shihang Wang. 2022b. <a href="#">Long time no see! open-domain conversation with long-term persona memory</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2639–2650, Dublin, Ireland. Association for Computational Linguistics.	889
838		890
839		891
840		892
841		893
842		894
843		895
844	Se-eun Yoon, Zhankui He, Jessica Echterhoff, and Julian McAuley. 2024. <a href="#">Evaluating large language models as generative user simulators for conversational recommendation</a> . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 1490–1504, Mexico City, Mexico. Association for Computational Linguistics.	896
845		
846		
847		
848		
849		
850		
851		
852		
853	Qiang Zhang, Jason Naradowsky, and Yusuke Miyao. 2023. <a href="#">Mind the gap between conversations for improved long-term dialogue generation</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 10735–10762, Singapore. Association for Computational Linguistics.	897
854		898
855		899
856		900
857		901
858		902
859	Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. <a href="#">Personalizing dialogue agents: I have a dog, do you have pets too?</a> In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2204–2213, Melbourne, Australia. Association for Computational Linguistics.	903
860		
861		
862		
863		
864		
865		
866		
867	Tenggan Zhang, Xinjie Zhang, Jinming Zhao, Li Zhou, and Qin Jin. 2024. <a href="#">ESCoT: Towards interpretable emotional support dialogue systems</a> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 13395–13412, Bangkok, Thailand. Association for Computational Linguistics.	
868		
869		
870		
871		
872		
873		
874	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. <a href="#">Bertscore: Evaluating text generation with BERT</a> . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	
875		
876		
877		
878		
879		
	Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth. 2019. <a href="#">“going on a vacation” takes longer than “going for a walk”</a> : A study of temporal commonsense understanding. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3363–3369, Hong Kong, China. Association for Computational Linguistics.	
	Pei Zhou, Karthik Gopalakrishnan, Behnam Hedayatnia, Seokhwan Kim, Jay Pujara, Xiang Ren, Yang Liu, and Dilek Hakkani-Tur. 2021. <a href="#">Commonsense-focused dialogues for response generation: An empirical study</a> . In <i>Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue</i> , pages 121–132, Singapore and Online. Association for Computational Linguistics.	
	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. <a href="#">SOTOPIA: Interactive evaluation for social intelligence in language agents</a> . In <i>The Twelfth International Conference on Learning Representations</i> .	

## A Data Construction Details

### A.1 ChatGPT Prompts

We provide prompts used for data construction processes of both evaluation and training datasets. The contents within curly brackets represent the corresponding elements for each example. The statements used for ATOMIC<sub>20</sub> duration estimation were constructed by concatenating the head and tail with a conjunction that represents each relation category. We present the prompts in the order of the processes, where the output of each step serves as the input for the next step.

#### ATOMIC<sub>20</sub> Duration Estimation

You are given a statement about common events in our daily lives. Your task is to estimate the typical duration of the key event in the form of (quantity of time + unit) (e.g., seconds, minutes, hours, days, weeks, months, years, decades, or centuries) based on the temporal common sense of average humans.

[Examples]

Statement: After dinner, he went to look for Max one last time before he had to take a bath and go to bed.

Key event: having dinner

Duration: 1 hour

Statement: Jennie and Bryan boarded a 6:00 A.M. flight from Seoul to Los Angeles International Airport.

Key event: flight from Seoul to Los Angeles

Duration: 12 hours

Event: Carl Laemmle, head of Universal Studios, gave Einstein a tour of his studio and introduced him to Chaplin.

Key event: tour of his studio

Duration: 45 minutes

[End of Examples]

Statement: {statement}

#### MC-TACO Event Descriptions

You are given an event and a question and answer for the duration that denotes how much time is needed for the event to

happen.

Write a story regarding the event in one sentence.

Sentence: {sentence}

Question: {question}

Answer: {duration}

#### ATOMIC<sub>20</sub> Event Descriptions

You are given a statement, the key event and the duration that denotes how much time is needed for the event to happen.

Write a story regarding the event in one sentence.

Statement: {statement}

Key event: {event}

Duration: {duration}

#### Dialog Generation

You are given an event narrative and the duration. Your task is to create an instant message dialog between two speakers. The following conditions MUST be met.

[Instructions]

1. Speaker {A,B} is in the middle of the event now, while speaker {B,A} is physically apart from.

2. Do not directly mention the duration in the dialog.

3. After {B,A}'s last turn, add "[{duration} later]", where duration is the amount of time passed in real world.

4-1. Generate {A,B}'s last message which is timely as if {A,B} spent time to finish the event.

4-2. In contrast, generate {A,B}'s last message as if {A,B} is responding instantaneously right before the event to happen.

Make sure that the timely response and the instantaneous response are time-situationally different.

[End of Instructions]

[Example]

{dialog example}

[End of Example]

Narrative: {event description}

Duration: {duration}

### A.2 Few-shot Examples

We provide six author-written dialog examples randomly fed into GPT-4o as one-shot demonstrations when generating dialogs for TIMELYCHAT using the MC-TACO dataset.

#### 5-turn Dialog

Narrative: After dinner, he took a shower before he went to bed.

Duration: 20 minutes

A: I finally got home. What a day!  
 B: It's eleven p.m. and you just got back home? It must be very tough day today.  
 A: Whooa Imma take a shower. I'm too tired.  
 B: Wash out all your fatigue with it.  
 [20 minutes later]  
 (delayed response)  
 A: I feel much better now! You didn't go to bed?  
 (instantaneous response)  
 A: How nice of you :) Give me a moment.  
 brb

#### 6-turn Dialog

Narrative: She has taken calculus class and she had a final exam.

Duration: 2 hours

A: Hey, what are you up to?  
 B: I'm gonna take the calculus final exam in 20 minutes. I feel so nervous.  
 A: You studied really hard, didn't you? I'm 100% sure you'll do well.  
 B: But the last two chapters were too difficult for me to understand.  
 A: That means others feel the same. Don't worry too much!  
 [2 hours later]  
 (delayed response)  
 B: It wasn't much harder than I expected. I hope I get a good grade.  
 (instantaneous response)  
 B: Thank you for cheering me up. I hope the exam is not that hard.

#### 7-turn Dialog

Narrative: He enjoyed working out at the gym.

Duration: 1 hour 30 minutes

A: I'm going to the gym now. Wanna join me?  
 B: I don't feel like working out today. Sorry.  
 A: You don't feel good? What happened?  
 B: I played football so hard yesterday that I can't even walk right.  
 A: Okay, I understand. Maybe next time!  
 B: Enjoy your routine! I think I can make it tomorrow.  
 [2 hours later]  
 (delayed response)  
 A: I focused on my lower body today. Chest tomorrow?  
 (instantaneous response)  
 A: Gonna work out hard on my lower body. Chest tomorrow?

#### 8-turn Dialog

Narrative: She had felt so tired that she went to bed right after the tv show.

Duration: 8 hours

A: Are you watching the saturday night live?  
 B: I'm watching it now but I'm too tired.  
 A: I didn't expect today's host is such a comedian lol  
 B: Yeah almost the end of the show. I feel like going to bed little bit early.  
 A: What made you so tired? You had any plan?  
 B: I went to an amusement park with my sister. We had a really good time there.  
 A: Oh I see. Think I should let you go. Sleep tight!  
 [8 hours later]  
 (delayed response)  
 B: Good morning. Did you sleep tight, too?  
 (instantaneous response)  
 B: Good night. I'll text you in the morning.

#### 9-turn Dialog

Narrative: He took an intercity bus to get to

his hometown.

Duration: 5 hours

A: What are you going to do on these holidays?

B: My parents and I usually have dinner together on the Eve.

A: Me too. So I'm heading to my town right now.

B: How do you get there? By bus or train?

A: I used to take trains, but I take an intercity bus for this time.

B: Why? the tickets' been already sold out?

A: Unfortunately yes... It will take little bit longer.

B: Have a nice trip though. Your family must be waiting for you.

[5 hours later]

(delayed response)

A: Finally I'm back at home! It took almost 5 hours.

(instantaneous response)

A: I'm gonna sleep all along in the bus. See you a few hours later.

### 10-turn Dialog

Narrative: She played the popular online

game with her friends.

Duration: 30 minutes

A: Have you heard of the League of Legends?

B: Absolutely! I play it almost everyday with my classmates.

A: I've heard of, but I've never played if before.

B: We have a game soon. Wanna join us?

A: Isn't it a team game? I'm not a good gamer though.

B: It's not a big deal. They will welcome you.

A: Well, maybe next time. I need to play it by myself first.

B: How about getting tutorial with me after this? I'll teach you.

A: Sounds good. Enjoy your game with your teammates.

[30 minutes later]

(delayed response)

B: We won! The game was nip and tuck. We were so close to losing.

(instantaneous response)

B: I'll be back just after the game. Wish me a good luck!

### A.3 Time Interval Distribution

We semantically categorize the durations and present the distributions of time intervals in the two proposed datasets in Table 6 for a comprehensive comparison.

Time Interval	Train	Test
0-5 min (instant)	18.5%	5.6%
5-30 min (short)	42.3%	38.6%
30 min-2 hrs (moderate)	23.8%	32.1%
2-6 hrs (long)	2.9%	13.9%
6-12 hrs (half-day)	2.0%	5.2%
12-24 hrs (full-day)	10.5%	4.6%

Table 6: Time interval distribution of the augmented training dataset (Train) and TimelyChat benchmark (Test).

## B Evaluation Setup Details

### B.1 Implementation Details

**Instruction-tuned Models.** We use vLLM library<sup>2</sup> for the inference of LLaMA 3.1 Instruct 8B and 70B on 4 NVIDIA A100 80GB GPUs. All instruction-tuned models employ top- $p$  sampling with temperature  $T = 1.0$  and  $p = 0.95$  during inference. We provide the prompts used for in-context learning methods on both response timing prediction and time-conditioned response generation below.

#### Prompt for Response Timing Prediction

You are given a conversation between two speakers.

Your task is to estimate a time interval needed until the next response, considering the duration of the event in the conversation ranging from 0 minutes to 24 hours (1 day). If the next response is expected to be immediate, you will output "0 minutes".

Otherwise, you will output a digit and a unit of time (e.g., 5 minutes, 2 hours).

Just output the time interval without any other text.

[Example  $n$ ]

{few-shot OR CoT example}

### Dialog context ###

{context}

Answer format:  $n$  ( $0 \leq n \leq 1440$ ) minutes

The estimated time interval is:

#### Prompt for Time-conditioned Response Generation

You are given a conversation between two speakers and the elapsed time since the last utterance.

Your task is to generate the next response that aligns well with the temporal context represented by the time interval in parentheses.

Just output the response without any other text.

[Example  $n$ ]

<sup>2</sup><https://docs.vllm.ai>

{few-shot OR CoT example}

### Dialog context ###

{context}

### Next response ###

{target speaker}: ({time interval} later)

**Time-aware Models.** We use Huggingface library<sup>3</sup> for the inference of MSC 3B and ReBot 400M<sup>4</sup>. We converted MSC 3B on the ParlAI framework<sup>5</sup> into a Huggingface checkpoint. We fine-tune GapChat 3B and TIMER 3B on each training data using the DeepSpeed library<sup>6</sup> with mixed precision training (Micikevicius et al., 2018). We train the models for 3 epochs using AdamW optimizer (Loshchilov and Hutter, 2019) with a learning rate of  $1e-4$ , running on 2 NVIDIA A100 80GB GPUs for 9 hours. We use  $\lambda = 1.0$  as a balanced scale factor of the two losses when training TIMER 3B. During inference, we apply beam search with the beam size of 3 and top- $p$  sampling with  $p = 0.95$ .

### B.2 User Simulator Prompts

We present the prompt fed into GPT-4o to create the user simulator.

#### User Simulator Prompt

You are a user simulator (user) engaging in an event-driven dialog with a dialog agent (agent).

Given the dialog context, your task is to proceed the conversation by one turn under the following assumptions:

1. agent responds after the elapsed time specified in the parentheses from the previous user utterance. If the delay is "0 minutes", agent is assumed to respond immediately.
2. user is assumed to respond to agent without any delay.

Conversation:

{context}

<sup>3</sup><https://huggingface.co>

<sup>4</sup><https://huggingface.co/jihyoung/rebot-generation>

<sup>5</sup><https://parl.ai>

<sup>6</sup><https://www.deepspeed.ai>

### B.3 LLM-as-a-Judge Details

We elucidate the LLM-as-a-judge prompts used in turn-level and dialog-level evaluations, along with the evaluation criteria and steps for each metric.

#### Turn-level Prompt

You will be given a conversation between two individuals via messaging, along with the elapsed time since the last utterance. You will then be given on potential response for the next turn.

Your task is to rate the response on one metric. Please make sure you read and understand these instructions carefully.

Evaluation Criteria:  
{metric} (1-5): {criteria}

Evaluation Steps:  
{steps}

#### Dialog-level Prompt

You will be given a conversation between a dialog agent and a user.

Throughout the conversation, the agent proactively determines the delay of its response to the user's previous message, simulating delayed responses due to event experiences that take certain time to process.

At each agent's turn, the delay is provided in the parentheses followed by the message. Your task is to rate the dialog agent on one metric. Please make sure you read and understand these instructions carefully.

Evaluation Criteria:  
{metric} (1-5): {criteria}

Evaluation Steps:  
{steps}

#### Evaluation Criteria and Steps

- **Naturalness** (1-5): the extent to which the response reads naturally given the dialog context.
  1. Assess the flow and coherence of the response in the conversation: Consider how seamlessly the response connects with the previous message.
  2. Evaluate the tone and style compatibility: Determine if the response's tone and style match

those of the previous messages.

3. Rate on a scale from 1 to 5, where 1 indicates the response is unnatural or inappropriate, and 5 indicates a perfectly natural continuation of the conversation.

- (Turn-level) **Time-specificity** (1-5): the extent to which the response ONLY makes sense when the specified time has passed, contrary to a time-agnostic response that makes sense regardless of time.
  1. Read the provided conversation and take note of the elapsed time since the previous message.
  2. Consider the context of the conversation, focusing on how the passage of time might affect the relevance or appropriateness of the response.
  3. Evaluate whether the potential response provided is time-specific. That is, determine if the response directly relates to or is clearly influenced by the elapsed time between the last utterance and the response.
  4. Rate on a scale from 1 to 5, where 1 indicates the response is completely time-agnostic and unaffected by the passage of time, and 5 indicates the response is entirely time-specific; it only makes sense because of the amount of time that has passed since the previous message.
- **Coherence** (1-5): the extent to which the agent maintains a good conversation flow.
  1. Assess the flow and coherence of the agent's responses in the conversation.
  2. Evaluate the tone and style compatibility throughout the conversation.
  3. Rate on a scale from 1 to 5, where 1 indicates the agent's responses are incoherent or inappropriate, and 5 indicates the agent's responses are perfectly coherent and appropriate.
- **Delay-appropriateness** (1-5): the extent to which the agent poses delays with appropriate frequency and amount.
  1. Assess whether the agent poses unnecessary or excessively frequent delays that could harm the conversation flow.
  2. Evaluate whether the amounts of delays (if not 0 minutes) reflect the typical duration of events implied in the corresponding message.
  3. Rate on a scale from 1 to 5, where 1 indicates the agent overuses and misuses delays, and 5 indicates the agent uses delays appropriately in terms of frequency and amount.
- (Dialog-level) **Time-specificity** (1-5): the extent to which the agent's responses ONLY make sense

1040 when the specified time has passed, contrary to a  
1041 time-agnostic responses that make sense regard-  
1042 less of time.

- 1043 1. Read the provided conversation and take note  
1044 of the elapsed times since the previous messages.
- 1045 2. Consider the context of the conversation, fo-  
1046 cusing on how the passage of time might affect  
1047 the relevance or appropriateness of the agent’s  
1048 responses.
- 1049 3. Evaluate whether the agent’s responses are  
1050 time-specific. That is, determine if the responses  
1051 directly relate to or are clearly influenced by the  
1052 elapsed times.
- 1053 4. Rate on a scale from 1 to 5, where 1 indi-  
1054 cates the agent’s responses are completely time-  
1055 agnostic and unaffected by the passage of time,  
1056 and 5 indicates the agent’s responses are entirely  
1057 time-specific; they only make sense because of  
1058 the amount of time that has passed since the pre-  
1059 vious message.

#### 1060 **B.4 Human Evaluation Forms**

1061 We present the questionnaire forms for human eval-  
1062 uation in Figure 4 and Figure 5. Note that the re-  
1063 sponses and agents are anonymized, which ensures  
1064 blind assessment.

### 1065 **C Case Study**

1066 Figure 6 presents illustrative examples of dialog  
1067 simulations conducted with TIMER 3B, LLaMA  
1068 3.1 70B, and GPT-4o for the same event. In TIMER  
1069 3B’s conversation, the agent correctly identifies  
1070 a situation where a delay is appropriate, specifi-  
1071 cally, when the user’s utterance (e.g., “...let me  
1072 know...”) suggests a natural pause in the conver-  
1073 sation. The agent then introduces a realistic 1-day  
1074 delay before responding with an update about find-  
1075 ing the missing person, successfully justifying the  
1076 delay. In contrast, LLaMA 3.1 70B generates de-  
1077 layed responses in every turn, but the predicted  
1078 time intervals appear somewhat arbitrary (e.g., 660  
1079 minutes, 77 minutes). Furthermore, its responses  
1080 lack time specificity, making it difficult to establish  
1081 a clear temporal correlation between the predicted  
1082 delays and the generated response. GPT-4o pre-  
1083 dicts more realistic time intervals that better align  
1084 with the temporal context compared to LLaMA 3.1  
1085 70B. However, it still fails to generate time-specific  
1086 responses, meaning the predicted delays are not  
1087 well justified. It also exhibits a tendency to overuse  
1088 delays, which can disrupt the natural flow of con-

versation. We observe similar behavior in LLaMA  
3.1 8B and GPT-3.5, reinforcing these findings.

1089  
1090

**Evaluation Methods and Guidelines**

Please read the conversation between the two speakers and choose the response you judge to be of higher quality based on the criteria provided below.

The conversation is conducted on a chatting platform (e.g., KakaoTalk, DM). Please assume that the two presented responses are made after the actual amount of time indicated in parentheses has elapsed since the previous utterance.

If both responses are equally excellent or both are poor, making it difficult to distinguish their quality, please select 'Tie'.

The evaluation consists of the following two metrics:

- **Naturalness**: Does the response fit the flow of the conversation and connect naturally with the preceding dialogue (regardless of the timing)? Are the tone and style consistent with the overall conversation?
- **Time-Specificity**: Is the response an utterance that can **ONLY** be made after the specified amount of time (e.g., X minutes later) has passed? Or is it a generic utterance that remains valid regardless of the time condition?

---

**Context**

A: How's the countryside life treating you?  
 B: It's been amazing! Joey and I are actually out running in the fields right now.  
 A: That sounds like fun. I miss the open spaces of the countryside.  
 B: You should definitely come visit. Running here is so freeing, with the wind and the wide-open spaces.  
 A: Sounds heavenly. Is Joey keeping up with you?  
 B: Barely! We're racing each other. It's a close match, but I think I've got the edge.  
 A: Just make sure you both don't run into a scarecrow or something.  
 B: Ha! I'll keep an eye out. This is the perfect way to enjoy the sunset.  
 A: Enjoy the rest of your run. Say hi to Joey for me!  
 (20 minutes later)  
 B: \_\_\_\_\_

---

Please select the superior response between the two options. \*

Response 1: We did it! Joey won the race, and we're both exhausted but happy. The sunset was breathtaking.  
 Response 2: Hey, we're back from our run! Joey says hi and we both agree that was the perfect way to end the day. The sky looked beautiful with the sunset. Wish you could've joined us!

---

	Response 1	Tie	Response 2
Naturalness	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Time-Specificity	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Figure 4: Questionnaire for turn-level human evaluation.

## Evaluation Methods and Guidelines

Two sets of dialogs between a **User (U:)** and an **Agent (A:)** will be presented. These two dialogs share the same conversation starter but diverge as they proceed with two different agents. Please note that the User responds without any time delay, while the Agent is assumed to respond with a delay only when a condition is specified in parentheses (e.g., (X minutes later)). In such cases, assume that the indicated amount of time has actually elapsed since the previous utterance.

Please evaluate which agent performed better based on the criteria below. If both agents are equally excellent or both are poor, making it difficult to distinguish their quality, please select 'Tie'.

The evaluation focuses solely on the **Agent's utterances** based on the following three metrics:

- **Coherence:** (Regardless of time conditions) Is the flow of the conversation natural? Are the tone and style of the dialog consistent?
- **Delay-Appropriateness:** Does the agent apply a time delay at the appropriate moment? Is the presented Time Delay reasonable and natural, considering the typical duration of the events occurring in the conversation?
- **Time-Specificity:** Are the agent's responses utterances that can **ONLY** be made after the specified amount of time (e.g., X minutes later) has passed? Or are they responses that would remain valid even if the time condition were absent or different?

### User vs. Agent Chats

[ Agent 1 vs. User Conversation ]

U: Hey, weren't you supposed to be at Jeff's place 10 minutes ago?  
A: Yeah, I was, but I got caught up at work. I'll be there soon.  
U: Can you let Jeff know you'll be a bit late or do you want me to tell him for you?  
A: I'll let him know. Thanks for understanding.  
U: No worries, just wanted to make sure you're covered. Let me know if you need anything else.  
A: (10 minutes later) Finally made it to Jeff's place. Sorry for the delay.  
U: Glad you made it! How's everything going over there?

[ Agent 2 vs. User Conversation ]

U: Hey, weren't you supposed to be at Jeff's place 10 minutes ago?  
A: Yeah, I got caught up with something, but I'm heading there right now. Thanks for the reminder!  
U: How far are you from his place? Do you need any directions?  
A: I'm about 15 minutes away, but I know the way, thanks! I'll let you know if I get lost.  
U: Sounds good. Hope you have a great time there! Let me know how it goes.  
A: (15 minutes later) Hey, I just got to Jeff's place! It wasn't too hard to find, thankfully. Thanks for checking in earlier!  
U: No problem! Enjoy the party, and say hi to everyone for me.  
A: Will do! I'll catch up with you later and let you know all about it. Have a good night!  
U: You too! Looking forward to hearing all about it. Take care!  
A: (240 minutes later) Hey! Just got back home. The party was awesome, and everyone says hi! We should definitely catch up soon; I've got some funny stories to share. Hope you had a good night too!  
U: I'm glad to hear you had a great time! I can't wait to hear all the stories. How does catching up over coffee tomorrow sound to you?

Please select the superior agent between the two. \*

Figure 5: Questionnaire for dialog-level human evaluation.

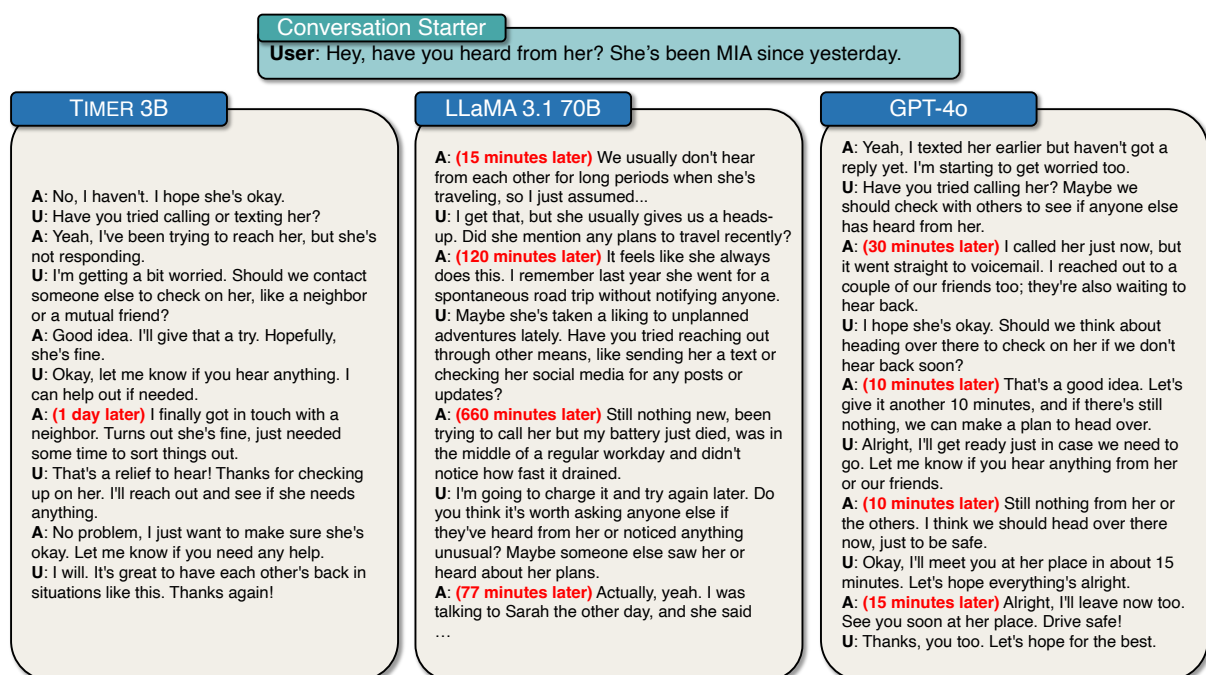


Figure 6: Examples of timely dialog simulations with GPT-4o as a user simulator. All examples begin with the same conversation starter and consist of 10 interactions. U and A represent the user and agent, respectively. We highlight the time intervals predicted by the agent in red, but omit them when the interval is 0 minutes.