

# Deep Think with Rehearsal for Low-Latency Team-AI Collaboration

Anonymous ACL submission

## Abstract

The integration of Large Language Models (LLMs) into scientific team meetings presents exciting opportunities to accelerate biomedical discovery, especially through their strong deep thinking capabilities enabled by multi-step reasoning and web search. However, such methods are computationally expensive and introduce substantial latency, limiting their effectiveness in real-time Team-AI communications. In this study, we propose *Deep Think with Rehearsal* (DTR), a novel framework that decouples deep reasoning from synchronous interaction in the AI4Science context. DTR transfers computationally intensive reasoning into an offline rehearsal phase, allowing the LLM to pre-cognize complex scientific contexts and deliver high-quality, "deep" insights with minimal latency during live interactions. To facilitate this research, we introduce the *Scientific Team Meeting Dataset* (STMD), a hybrid benchmark comprising authentic transcripts from three real-world biomedical research labs alongside extensive simulated multi-party deliberations synthesized from PubMed literature. Experiments in both simulated and real-world settings demonstrate that DTR consistently improves response quality while reducing inference latency compared to state-of-the-art methods, highlighting the effectiveness of rehearsal in enabling low-latency, high-quality scientific collaboration.<sup>1</sup>

## 1 Introduction

Large Language Models (LLMs) have demonstrated strong reasoning and tool augmentation capabilities, such as chain-of-thought problem-solving (Wei et al., 2022; Zhang et al., 2025) and tool use via web search (Yao et al., 2022; Schick et al., 2023), making them increasingly effective at enabling Team-AI collaboration and accelerating scientific discovery (Zheng et al., 2025; Singh et al.,

<sup>1</sup>Code and data are provided in the supplementary material and will be made publicly available upon acceptance.

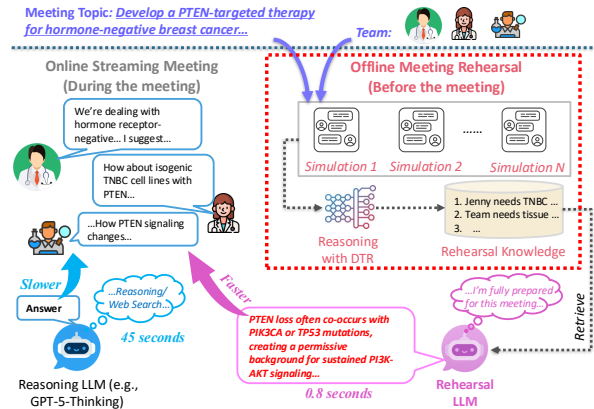


Figure 1: Comparison of traditional reasoning LLM and rehearsal LLM in real-time Team-AI collaboration. The traditional reasoning LLM (left) performs deep thinking online (during the meeting) with high latency, whereas the rehearsal LLM (right) rehearses offline (before the meeting) via meeting simulations and construct a rehearsal knowledge database, enabling fast retrieval during online interaction.

2025; Qi et al., 2024). However, integrating these capabilities into real-time conversations (e.g., scientific meetings) remains challenging, as deep reasoning and web search at inference time are computationally expensive, introducing prohibitive delays that disrupt conversational flow and limit their effectiveness in Team-AI collaboration (Gao et al., 2025; Bi et al., 2024). This constraint highlights a gap between how current LLM assistants behave and how effective human collaborators participate in meetings (Schmutz et al., 2024; Liao et al., 2023). Crucially, unlike general-purpose open-domain dialogues, scientific meetings are characterized by a high degree of predictability, as they are typically centered around predefined research objectives, specific datasets, or persistent scientific hypotheses. Once the meeting topic is known, human experts typically leverage this thematic stability to prepare in advance by surveying relevant literature, organizing key evidence, and formulating candi-

061 date hypotheses (Maslych et al., 2025; Král et al.,  
062 2023). This pre-meeting preparation ensures they  
063 can engage efficiently during live discussions rather  
064 than interrupting the flow to search for information.  
065 We argue that LLM-based assistants should follow  
066 a similar "preparation-first" principle. By exploit-  
067 ing the predictable scope of scientific deliberations,  
068 we can transfer the computational burden of deep  
069 reasoning into an offline phase, achieving high-  
070 quality insights while maintaining the low-latency  
071 responsiveness required for seamless real-time col-  
072 laboration (Pelikan and Hofstetter, 2023; Liu et al.,  
073 2025).

074 Motivated by this insight, we introduce *Deep*  
075 *Think with Rehearsal* (DTR), a rehearsal-based  
076 framework that transfers expensive online reason-  
077 ing and knowledge gathering to an offline rehearsal  
078 phase, enabling fast responses during deployment  
079 while preserving reasoning quality (Prince et al.,  
080 2024). As illustrated in Figure 1, in the **online**  
081 setting (i.e., real-time interaction during the sci-  
082 entific meeting), a conventional reasoning LLM  
083 (e.g., GPT-5-Thinking (OpenAI, 2025c) or Gemini  
084 3 (Google DeepMind, 2025)) performs deep rea-  
085 soning and web search during the discussion, which  
086 can take minutes per turn and introduce substan-  
087 tial response delays in Team-AI interaction (Cheng  
088 et al., 2025). In contrast, DTR shifts these compu-  
089 tationally intensive steps to an **offline** phase (i.e.,  
090 pre-meeting preparation), where it runs confidence-  
091 based meeting simulations conditioned on the meet-  
092 ing topic and participants and uses a reasoning  
093 LLM to perform deep reasoning and web search  
094 over simulated trajectories, thereby constructing  
095 a rehearsal knowledge database before the meet-  
096 ing starts. During the meeting, the DTR LLM  
097 only needs to leverage lightweight retrieval from  
098 the prepared knowledge database to deliver timely,  
099 high-quality suggestions, for example, recommend-  
100 ing that the team consider PIK3CA or TP53 when  
101 designing PTEN-targeted therapy studies for hor-  
102 mone receptor negative breast cancer, with a re-  
103 sponse latency of only 0.8 seconds.

104 To support the study of low-latency deep think-  
105 ing, we introduce the *Scientific Team Meeting*  
106 *Dataset* (STMD), a hybrid benchmark comprising  
107 authentic transcripts from three real-world biomed-  
108 ical research labs covering cancer, sepsis, and  
109 cardio-metabolic disease, along with extensive sim-  
110 ulated multi-party deliberations synthesized from  
111 PubMed literature (Sayers et al., 2024). STMD

112 features multi-role biomedical team meetings and  
113 is designed to evaluate both response quality and  
114 inference latency. In both simulated and real-world  
115 studies, DTR consistently improves reasoning qual-  
116 ity while substantially reducing inference latency  
117 compared with state-of-the-art methods (OpenAI,  
118 2025c), highlighting offline rehearsal as an effec-  
119 tive approach for real-time, high-quality Team-AI  
120 scientific collaboration.

121 Our contributions in this paper can be summa-  
122 rized as follows:

123 • **Conceptual Paradigm:** We formalize the  
124 problem of low-latency deep thinking in real-  
125 time collaboration and propose the "Rehearse-then-  
126 Interact" paradigm. This shifts the focus from op-  
127 timizing online inference to exploiting the inher-  
128 ent predictability of structured scientific meetings  
129 through offline preparation.

130 • **Framework:** We propose DTR, a novel frame-  
131 work that operationalizes this paradigm. By de-  
132 coupling expensive multi-step reasoning and web-  
133 based knowledge gathering into an offline phase,  
134 DTR enables LLMs to deliver "deep" insights with  
135 sub-second latency during live interactions.

136 • **Benchmark:** We introduce the STMD, a spe-  
137 cialized benchmark derived from PubMed litera-  
138 ture and real-world lab meetings. It provides a  
139 robust testbed for evaluating both the reasoning  
140 depth and temporal efficiency of LLMs in multi-  
141 party scientific deliberations.

142 • **Empirical Validation:** We conduct extensive  
143 experiments in both simulated environments and  
144 real-world biomedical lab settings. Results demon-  
145 strate that DTR consistently outperforms state-of-  
146 the-art methods in both response quality and infer-  
147 ence efficiency, proving the viability of rehearsal-  
148 based AI collaboration.

## 149 2 Related Work

150 **LLMs with Deep Think** Recent work improves  
151 LLM reasoning through structured intermediate  
152 steps and external knowledge integration. Chain-  
153 of-thought prompting (Wei et al., 2022), self-  
154 consistency (Wang et al., 2022), and Tree of  
155 Thoughts (Yao et al., 2023) elicit multi-step rea-  
156 soning and explore alternative paths via search.  
157 For grounding and tool use, retrieval-augmented  
158 generation (Lewis et al., 2020) and ReAct (Yao  
159 et al., 2022) integrate retrieval and tool-mediated  
160 actions into generation, while Toolformer (Schick  
161 et al., 2023), MRKL (Karpas et al., 2022), and

program-aided models (Gao et al., 2023) enable self-supervised tool use, modular tool composition, and code-based reasoning. Confidence-aware methods (Subramani et al., 2025; Taubenfeld et al., 2025) adjust reasoning depth to uncertainty to reduce low-utility computation, but over-reasoning can still inflate cost (Lacombe et al., 2025). These pipelines typically require iterative calls, retrieval, and tool execution during online inference, increasing latency and disrupting real-time interaction, whereas DTR shifts these costs to offline rehearsal and answers online via lightweight retrieval.

**LLMs with Rehearsal** Rehearsal is widely used in reinforcement learning and agent systems to simulate future interactions and improve decisions. In model-based RL, MuZero (Schrittwieser et al., 2020) and Dreamer (Hafner et al., 2023) plan via imagined rollouts, and Language Agent Tree Search (Zhou et al., 2023) performs lookahead over action branches. Beyond RL, offline simulation and self-play have been used to anticipate user responses and build reusable strategies in meeting settings (He et al., 2024, 2025; Kim et al., 2025). Role-play and continual exploration further accumulate interaction experience in multi-agent systems (Li et al., 2023; Park et al., 2023; Wang et al., 2023a; Yamada et al., 2025b), while offline training can strengthen multi-step reasoning (Wang et al., 2025). These lines of work often keep substantial reasoning on the online path or prioritize capability over latency, whereas DTR distills evidence-backed reasoning offline and answers online via lightweight retrieval.

### 3 Methodology

#### 3.1 Preliminary Definitions

We consider a fixed team  $\mathcal{T}$  with a set of roles (team members)  $\mathcal{R} = \{r_1, \dots, r_{|\mathcal{R}|}\}$ . A multi-role scientific meeting is a sequence of utterances  $\mathcal{D} = \{u_1, \dots, u_T\}$ , where each utterance is a speaker-content pair  $u_t = (r_t, x_t)$ ,  $r_t \in \mathcal{R}$  and  $x_t$  is a text utterance. Let  $h_t = \{u_1, \dots, u_t\}$  denote the meeting context up to round  $t$ .

**Task Definition.** Given the current *online* meeting context<sup>2</sup>  $h_t^{\text{on}}$  up to round  $t$  and a user query  $q_t$  issued at round  $t$ , our task is to generate an assistant response  $y_t$  with high quality under a strict

<sup>2</sup>We use the superscript “on” to indicate online meeting context (e.g.,  $h_t^{\text{on}}$ ); symbol without this superscript (e.g.,  $h_t$ ) refer to historical or offline rehearsal meeting context.

low-latency constraint for real-time interaction. We denote the deployed response LLM as  $\Omega$ , which takes the meeting context  $h_t^{\text{on}}$  and the query  $q_t$  as input and outputs the response:

$$y_t = \Omega(h_t^{\text{on}}, q_t). \quad (1)$$

Our objective is to learn an  $\Omega$  with high response quality while maintaining fast online inference.

**Confidence Score.** Given a LLM with parameters  $\theta$ , we quantify the confidence of an utterance  $u_t = (r_t, x_t)$  by the token-averaged log-likelihood of its content  $x_t = \{x_{t,1}, \dots, x_{t,|x_t|}\}$  under the meeting context  $(h_{t-1}, r_t)$ , where  $x_{t,<i} = \{x_{t,1}, \dots, x_{t,i-1}\}$ . Inspired by Fu et al. (2025), we define the utterance confidence score as:

$$C_\theta(u_t) = \frac{1}{|x_t|} \sum_{i=1}^{|x_t|} \log p_\theta(x_{t,i} | h_{t-1}, r_t, x_{t,<i}), \quad (2)$$

where  $p_\theta(\cdot)$  denotes the next-token probability induced by the LLM, and  $|x_t|$  is the number of tokens in the utterance. Larger  $C_\theta(u_t)$  indicates higher confidence, i.e., the model assigns higher probability to the utterance content.

**Team Memory and Persona.** We assume access to a historical meeting set for the same team  $\mathcal{T}$ ,  $\mathcal{H}_\mathcal{T} = \{\mathcal{D}^{(1)}, \dots, \mathcal{D}^{(N)}\}$ . From  $\mathcal{H}_\mathcal{T}$ , we build (i) a team memory bank  $\mathcal{M}_\mathcal{T}$  and (ii) role-specific persona banks  $\{\mathcal{P}_r\}_{r \in \mathcal{R}}$ , each with long-term and short-term components:

$$\begin{aligned} \mathcal{M}_\mathcal{T} &= \mathcal{M}_\mathcal{T}^L \cup \mathcal{M}_\mathcal{T}^S, \\ \mathcal{P}_r &= \mathcal{P}_r^L \cup \mathcal{P}_r^S. \end{aligned} \quad (3)$$

$\mathcal{M}_\mathcal{T}^L$  stores stable team knowledge (e.g., recurring hypotheses, validated protocols, shared constraints), while  $\mathcal{M}_\mathcal{T}^S$  stores recent updates (e.g., newest experimental results). Likewise,  $\mathcal{P}_r^L$  captures long-standing role expertise and preferences, while  $\mathcal{P}_r^S$  captures short-term focuses and temporary responsibilities. Given context  $h_t$  and role  $r$ , a retrieval module returns relevant memory and persona items:

$$\begin{aligned} m_t(r) &= \mathcal{R}_{\text{mem}}(\mathcal{M}_\mathcal{T}, h_t, r), \\ p_t(r) &= \mathcal{R}_{\text{per}}(\mathcal{P}_r, h_t, r). \end{aligned} \quad (4)$$

where  $m_t(r)$  and  $p_t(r)$  denote retrieved items from  $\mathcal{M}_\mathcal{T}$  and  $\mathcal{P}_r$ . We implement  $\mathcal{R}_{\text{mem}}(\cdot)$  and  $\mathcal{R}_{\text{per}}(\cdot)$  as embedding-based similarity retrieval (Xiao et al., 2023) over the long-term and short-term banks with temporal filters.

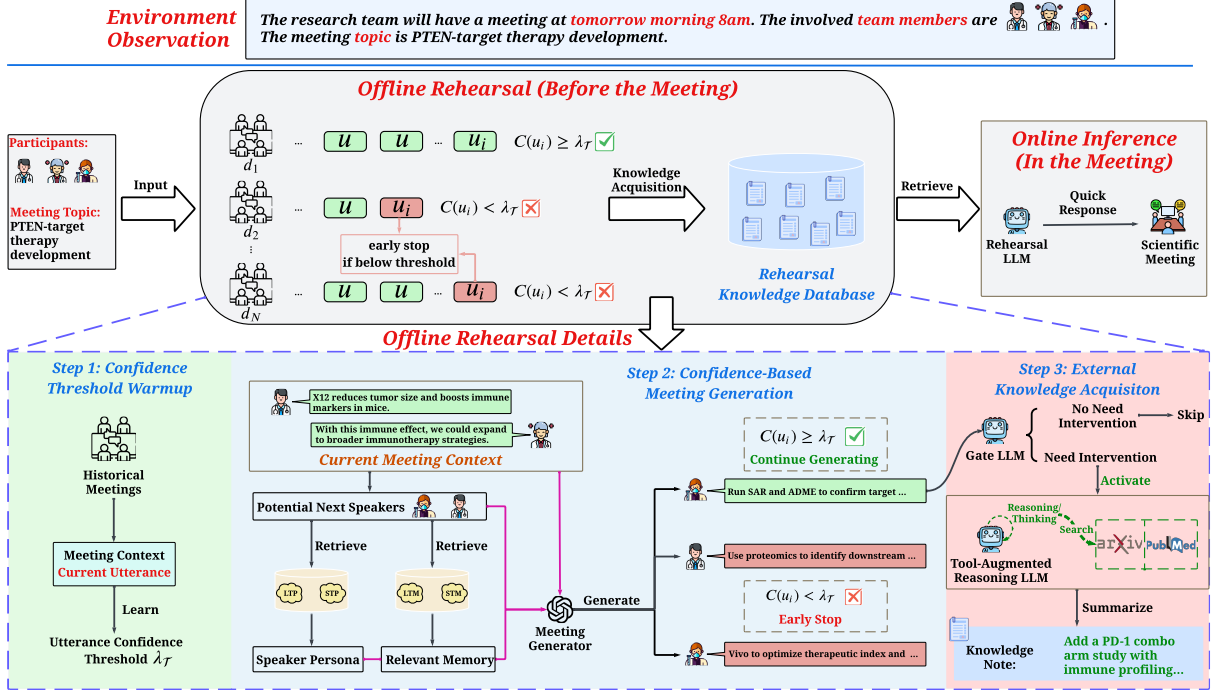


Figure 2: Illustration of our proposed DTR framework. Offline rehearsal stage: DTR simulates multi-role meetings with confidence, selectively invokes tool-augmented reasoning via a Gate LLM, and stores external knowledge notes into a rehearsal knowledge database. Online meeting stage: The rehearsal LLM conducts lightweight retrieval to gather relevant knowledge and generate low-latency, high-quality responses.

## 3.2 Deep Think with Rehearsal

In this section, we provide an overview of DTR, as illustrated in Figure 2. DTR shifts online reasoning and web search to an *offline rehearsal* stage, distills and stores the acquired knowledge, and reuses it for *fast online inference*. Section 3.2.1 covers three steps: (1) warming up team-specific confidence thresholds from historical meetings, (2) generating high-confidence rehearsal meetings conditioned on retrieved team memory and persona, and (3) selectively triggering a reasoning LLM to collect external knowledge and store distilled notes in a rehearsal database. Section 3.2.2 retrieves these notes and conditions the response model to produce low-latency, high-quality online meeting responses.

### 3.2.1 Offline Thinking with Rehearsal

**Confidence Threshold Warmup.** We first compute a *team-specific* confidence threshold from historical meetings. For each historical meeting  $\mathcal{D}^{(n)} \in \mathcal{H}_{\mathcal{T}}$ , we compute the utterance confidence  $C_{\theta}(u)$  for every utterance  $u \in \mathcal{D}^{(n)}$  using Equation 2. We then define the team threshold  $\lambda_{\mathcal{T}}$  as the empirical mean confidence over all historical

utterances:

$$\lambda_{\mathcal{T}} = \frac{1}{\sum_{n=1}^N |\mathcal{D}^{(n)}|} \sum_{n=1}^N \sum_{u \in \mathcal{D}^{(n)}} C_{\theta}(u), \quad (5)$$

where  $|\mathcal{D}^{(n)}|$  denotes the number of utterances in meeting  $\mathcal{D}^{(n)}$ . This warmup yields a simple but effective rule: a newly generated utterance  $\hat{u}$  is accepted only if its confidence exceeds the team’s historical average, i.e.,  $C_{\theta}(\hat{u}) \geq \lambda_{\mathcal{T}}$ .

**Meeting Generation with Confidence.** Given a predefined meeting topic and a team’s project mini proposal (e.g., goal, background, and available data), we initialize the rehearsal context  $h_0$  by formatting them as a kickoff context, consisting of a brief topic statement followed by the mini proposal. We then generate rehearsal meetings offline via tree-structured expansion. At each step  $t$ , for every candidate role  $r \in \mathcal{R}$ , we retrieve context-relevant team memory and persona items (Equation 4), denoted as  $m_t(r)$  and  $p_t(r)$ , and invoke a meeting generator  $G_{\phi}$  to propose the next utterance:

$$\begin{aligned} \hat{x}_{t+1}^{(r)} &= G_{\phi}(h_t, r, m_t(r), p_t(r)), \\ \hat{u}_{t+1}^{(r)} &= (r, \hat{x}_{t+1}^{(r)}). \end{aligned} \quad (6)$$

We compute the confidence  $C_\theta(\hat{u}_{t+1}^{(r)})$  for each utterance and apply a thresholded expansion rule:

$$h_{t+1} = \begin{cases} h_t \cup \{\hat{u}_{t+1}^{(r)}\}, & \text{if } C_\theta(\hat{u}_{t+1}^{(r)}) \geq \lambda_{\mathcal{T}}, \\ \text{terminate branch,} & \text{otherwise.} \end{cases} \quad (7)$$

where the accepted utterance is appended to the context to create child nodes, while rejected ones terminate their branches early. This prevents low-confidence utterances from accumulating and destabilizing long-horizon multi-role simulation. We repeat expansion until a maximum length  $T_{\max}$  (default  $T_{\max} = 12$ ), producing a set of high-confidence, team-conditioned rehearsal meetings  $\tilde{\mathcal{H}}_{\mathcal{T}}$  for subsequent knowledge acquisition.

### Rehearsal Knowledge Database Construction.

To avoid invoking expensive deep reasoning for every generated utterance, we follow the proactive assistant paradigm (Lu et al., 2024; Zhang et al., 2024) and introduce an *intervention gate* that decides when to call a reasoning LLM. Let  $h_t$  denote the rehearsal context at step  $t$ . We implement the gate as a small (0.5B) LLM  $\mathcal{Q}$  that predicts whether deeper reasoning is needed:

$$z_t = \mathcal{Q}(h_t), \quad z_t \in \{0, 1\}, \quad (8)$$

where  $z_t = 1$  triggers a reasoning call and  $z_t = 0$  skips it. When  $z_t = 1$ , we query a tool-augmented reasoning LLM  $\mathcal{G}$  that performs web search and multi-step reasoning to generate an intervention response:

$$\tilde{y}_t = \mathcal{G}(h_t), \quad (9)$$

where  $\tilde{y}_t$  includes the tool-mediated reasoning content and the final intervention message. We distill  $\tilde{y}_t$  into a compact knowledge note  $s_t$  and add it to the rehearsal knowledge database  $\mathcal{K}_{\mathcal{T}}$ :

$$\mathcal{K}_{\mathcal{T}} \leftarrow \mathcal{K}_{\mathcal{T}} \cup \{s_t\}. \quad (10)$$

$\mathcal{K}_{\mathcal{T}}$  stores offline-distilled notes from rehearsal for efficient online knowledge acquisition.

### 3.2.2 Online Inference with Rehearsal Knowledge

In the online meeting phase, we generate high-quality responses with low inference latency by avoiding web search and deep reasoning at inference time. Instead, we retrieve relevant rehearsal knowledge notes from  $\mathcal{K}_{\mathcal{T}}$  using lightweight embedding-based retrieval (Xiao et al., 2023) and

use them as in-context guidance for response generation. Given the online meeting context  $h_t^{\text{on}}$  and user query  $q_t$ , we form the retrieval query

$$\xi_t = (h_t^{\text{on}}, q_t). \quad (11)$$

and retrieve the top- $K$  notes:

$$\mathcal{S}_t = \{s^{(1)}, \dots, s^{(K)}\} = \mathcal{R}_{\text{kn}}(\mathcal{K}_{\mathcal{T}}, \xi_t), \quad (12)$$

where  $\mathcal{R}_{\text{kn}}(\cdot)$  is an embedding-based similarity retriever over  $\mathcal{K}_{\mathcal{T}}$ . We then generate the response by conditioning the deployed response LLM  $\Omega$  on the meeting context, the user query, and the retrieved notes:

$$y_t = \Omega(h_t^{\text{on}}, q_t, \mathcal{S}_t). \quad (13)$$

This design keeps online inference lightweight while leveraging offline-distilled external knowledge for low-latency response generation.

## 4 Dataset Construction

LLMs have demonstrated strong capability for synthesizing high-quality role-play datasets (Tao et al., 2024; Lim et al., 2024). However, existing biomedical meeting corpora such as MedDialog (Zeng et al., 2020), MediTOD (Saley et al., 2024), and MaLP (Zhang et al., 2023) primarily focus on doctor-patient interactions, while team-based scientific meetings remain largely underexplored. To fill this gap, we construct the *Scientific Team Meeting Dataset* (STMD), which combines a PubMed-grounded synthetic subset (STMD-Sim) with a real-world subset (STMD-Real) collected from three biomedical research labs. For STMD-Sim, we curate two variants, STMD<sub>1</sub> and STMD<sub>2</sub>, generated with GPT-5 and GPT-4o, respectively.

For STMD-Sim, each team project is anchored to one PubMed paper (Sayers et al., 2024), as illustrated in Figure 3. For each paper, we use a Prophet LLM (with access to the full paper) to extract three paper-grounded artifacts: a compact mini-proposal describing the project goal, essential background, and relevant data/resources; a golden conclusion summarizing the key findings and results; and an ordered sequence of meeting topics that guides a discussion toward the project goal. Based on these artifacts, we simulate meetings with four domain-specific roles: Pharmacologist, Medicinal Chemist, Bioinformatician, and Clinical Physician. Given the current topic and meeting context, the Prophet LLM selects the next speaker, and a Role-Play LLM (without access to the paper) generates the

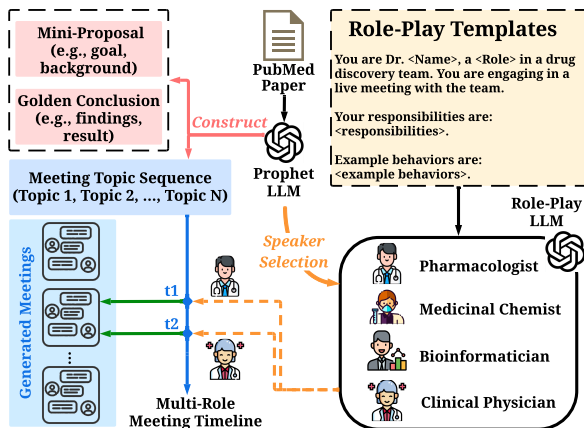


Figure 3: Overview of STMD-Sim generation. For each PubMed paper, a Prophet LLM constructs paper-grounded project artifacts (mini-proposal, golden conclusion, and a meeting-topic sequence) and performs speaker selection at each turn. A Role-Play LLM, guided by role-specific templates, then generates the corresponding utterance for selected roles, producing multi-role meetings along the timeline.

corresponding utterance using role-specific templates. To match our rehearsal setting, we generate ten offline meetings per project as historical meetings for confidence-threshold warmup and memory or persona construction, and one additional online meeting as the held-out evaluation meeting. To assess the synthetic dataset quality, we recruit two biomedical PhD researchers<sup>3</sup> to independently rate 30 randomly sampled meetings from each STMD-Sim variant. They use the rubric in Table 6 to score Meeting Coherence, Factual Faithfulness, and Role Consistency on a 1–5 scale. Table 1 shows that the STMD-Sim variants receive high quality scores on STMD<sub>1</sub>/STMD<sub>2</sub> for Coherence (4.27/4.58), Faithfulness (4.65/4.73), and Role Consistency (4.73/4.73). We also report inter-annotator agreement using Cohen’s  $\kappa$ , obtaining 0.35/0.38 (Coherence), 0.69/0.54 (Faithfulness), and 0.69/0.85 (Role Consistency) on STMD<sub>1</sub>/STMD<sub>2</sub>.

In addition to STMD-Sim, we collect STMD-Real from three biomedical research labs. STMD-Real contains multi-party lab meeting transcripts collected over 12 weeks from three research teams covering cancer, sepsis, and cardio-metabolic disease, where each team includes one Principal Investigator (PI) and five PhD student researchers<sup>4</sup>. We treat meetings from the first 10 weeks as historical

<sup>3</sup>More details of the human evaluation protocol for data quality are provided in Appendix B.1.

<sup>4</sup>Anonymization and participant consent are provided in Ethics Statement.

Dataset Quality Assessment (Scale 1-5)			
Metric	STMD <sub>1</sub>	STMD <sub>2</sub>	Average
Meeting Coherence	4.27	4.58	4.42
Factual Faithfulness	4.65	4.73	4.69
Role Consistency	4.73	4.73	4.73
Inter-Annotator Agreement (Cohen’s $\kappa$ )			
Metric	STMD <sub>1</sub>	STMD <sub>2</sub>	Average
Meeting Coherence	0.35	0.38	0.37
Factual Faithfulness	0.69	0.54	0.62
Role Consistency	0.69	0.85	0.77

Table 1: Human evaluation of dataset quality for two STMD-Sim variants. Two annotators rate meeting-level quality on a 5-point scale, and inter-annotator agreement is measured by Cohen’s  $\kappa$ .

and those from the remaining 2 weeks as online evaluation. Detailed dataset statistics for STMD-Sim and STMD-Real are reported in Appendix A.

## 5 Experiment

### 5.1 Experimental Settings

In offline rehearsal, we generate 12 multi-round team meetings per project. For retrieval-augmented generation, we use a default retrieval size of 3 across all methods. We compare our method with baselines on three LLM backbones: Llama 3.1-8B-Instruct (Meta, 2024), Gemma 3-12B-it (Google, 2025), and gpt-oss-120b (OpenAI, 2025b). Unless otherwise specified, all methods share the same external web-search interface, PubMed MCP (cyanheads, 2025). Because our source data span January 1, 2024 to January 1, 2025, we restrict web search to publications dated before January 1, 2024 to avoid temporal leakage. We set temperature to 1.0 and top\_p to 1.0, enable streaming generation, and accelerate inference with vLLM (Kwon et al., 2023). All experiments are implemented with Transformers (Wolf et al., 2020) and PyTorch (Paszke et al., 2017), and run on four 80GB NVIDIA A100 GPUs.

### 5.2 Baselines.

We compare DTR with several representative baselines: (i) Standard, a retrieve-then-read pipeline that issues a single search using query, concatenates the top-3 retrieved snippets as evidence, and generates the final response conditioned on this evidence; (ii) ReAct (Yao et al., 2022), which interleaves reasoning and tool use by deciding whether to issue additional searches based on intermediate reasoning and observed evidence, then refining the answer accordingly; (iii) SAS (Self-Ask

LLM Backbone	Method	Response Quality (%) $\uparrow$								Response Latency (Sec.) $\downarrow$			
		ROUGE-1		ROUGE-L		BLEU-1		BLEU-N		FTL		E2EL	
		STMD <sub>1</sub>	STMD <sub>2</sub>	STMD <sub>1</sub>	STMD <sub>2</sub>	STMD <sub>1</sub>	STMD <sub>2</sub>	STMD <sub>1</sub>	STMD <sub>2</sub>	STMD <sub>1</sub>	STMD <sub>2</sub>	STMD <sub>1</sub>	STMD <sub>2</sub>
gpt-oss-120b	Standard	41.8 $\pm$ 0.2	19.5 $\pm$ 0.5	29.4 $\pm$ 0.2	13.2 $\pm$ 0.2	7.4 $\pm$ 0.2	17.0 $\pm$ 0.2	6.5 $\pm$ 0.2	15.6 $\pm$ 0.0	9.7 $\pm$ 0.0	8.1 $\pm$ 0.0	11.1 $\pm$ 0.5	9.5 $\pm$ 0.2
	ReAct	46.8 $\pm$ 0.1	20.7 $\pm$ 0.7	35.3 $\pm$ 0.1	13.7 $\pm$ 0.0	7.8 $\pm$ 0.2	17.6 $\pm$ 0.3	7.1 $\pm$ 0.2	15.9 $\pm$ 0.0	19.6 $\pm$ 0.2	22.3 $\pm$ 0.0	21.4 $\pm$ 0.0	23.7 $\pm$ 0.1
	SAS	48.0 $\pm$ 0.2	21.6 $\pm$ 0.6	<u>36.9 <math>\pm</math>0.2</u>	13.7 $\pm$ 0.3	8.0 $\pm$ 0.1	16.8 $\pm$ 0.2	7.4 $\pm$ 0.1	15.7 $\pm$ 0.0	22.4 $\pm$ 0.0	24.2 $\pm$ 0.1	23.8 $\pm$ 0.0	25.8 $\pm$ 0.1
	DRA	<u>48.6 <math>\pm</math>0.2</u>	<u>24.4 <math>\pm</math>0.5</u>	35.2 $\pm$ 0.2	<u>15.0 <math>\pm</math>0.2</u>	<u>8.3 <math>\pm</math>0.2</u>	<u>19.4 <math>\pm</math>0.2</u>	<u>7.8 <math>\pm</math>0.2</u>	<u>17.5 <math>\pm</math>0.1</u>	24.2 $\pm$ 0.1	24.7 $\pm$ 0.0	26.2 $\pm$ 0.3	26.4 $\pm$ 0.1
	DTR (Ours)	<b>52.1 <math>\pm</math>0.4</b>	<b>26.8 <math>\pm</math>0.3</b>	<b>37.8 <math>\pm</math>0.1</b>	<b>16.5 <math>\pm</math>0.2</b>	<b>12.5 <math>\pm</math>0.5</b>	<b>21.3 <math>\pm</math>0.1</b>	<b>11.2 <math>\pm</math>0.4</b>	<b>19.3 <math>\pm</math>0.1</b>	<b>1.0 <math>\pm</math>0.0</b>	<b>1.1 <math>\pm</math>0.1</b>	<b>1.9 <math>\pm</math>0.1</b>	<b>2.1 <math>\pm</math>0.0</b>
Gemma-3-12b-it	Standard	42.1 $\pm$ 0.1	22.3 $\pm$ 0.9	31.1 $\pm$ 0.1	14.7 $\pm$ 0.6	15.9 $\pm$ 0.8	33.3 $\pm$ 0.8	14.0 $\pm$ 0.5	24.0 $\pm$ 0.1	<u>10.2 <math>\pm</math>0.2</u>	<u>8.9 <math>\pm</math>0.1</u>	<u>10.7 <math>\pm</math>0.2</u>	<u>9.6 <math>\pm</math>0.2</u>
	ReAct	<u>51.2 <math>\pm</math>0.2</u>	26.6 $\pm$ 0.1	<u>42.4 <math>\pm</math>0.2</u>	19.1 $\pm$ 0.6	16.6 $\pm$ 0.7	34.6 $\pm$ 0.1	15.4 $\pm$ 0.7	26.7 $\pm$ 0.3	19.2 $\pm$ 0.1	24.8 $\pm$ 0.2	20.1 $\pm$ 0.1	26.2 $\pm$ 0.7
	SAS	45.5 $\pm$ 0.1	25.3 $\pm$ 0.6	37.1 $\pm$ 0.1	17.2 $\pm$ 0.3	18.8 $\pm$ 0.6	35.7 $\pm$ 0.5	16.9 $\pm$ 0.4	25.6 $\pm$ 0.3	22.0 $\pm$ 0.2	22.8 $\pm$ 0.0	23.4 $\pm$ 0.1	24.1 $\pm$ 0.5
	DRA	50.8 $\pm$ 0.1	<u>34.2 <math>\pm</math>0.4</u>	41.0 $\pm$ 0.1	<u>23.5 <math>\pm</math>0.3</u>	<u>19.6 <math>\pm</math>0.2</u>	<b>37.9 <math>\pm</math>0.4</b>	<u>17.4 <math>\pm</math>0.1</u>	<b>27.8 <math>\pm</math>0.1</b>	29.1 $\pm$ 0.0	31.2 $\pm$ 0.0	30.8 $\pm$ 0.6	32.8 $\pm$ 0.2
	DTR (Ours)	<b>54.2 <math>\pm</math>0.3</b>	<b>36.8 <math>\pm</math>0.1</b>	<b>44.5 <math>\pm</math>0.2</b>	<b>25.4 <math>\pm</math>0.1</b>	<b>22.8 <math>\pm</math>0.4</b>	<u>36.2 <math>\pm</math>0.2</u>	<b>19.7 <math>\pm</math>0.5</b>	<u>27.0 <math>\pm</math>0.1</u>	<b>1.1 <math>\pm</math>0.2</b>	<b>1.2 <math>\pm</math>0.0</b>	<b>1.9 <math>\pm</math>0.1</b>	<b>2.0 <math>\pm</math>0.1</b>
Llama-3.1-8B-Instruct	Standard	38.8 $\pm$ 0.2	29.8 $\pm$ 0.6	31.6 $\pm$ 0.2	25.5 $\pm$ 0.7	10.8 $\pm$ 0.8	20.0 $\pm$ 0.7	9.3 $\pm$ 0.7	18.5 $\pm$ 0.3	<u>7.6 <math>\pm</math>0.2</u>	<u>8.1 <math>\pm</math>0.0</u>	<u>8.2 <math>\pm</math>0.2</u>	<u>8.9 <math>\pm</math>0.2</u>
	ReAct	43.1 $\pm$ 0.2	31.0 $\pm$ 0.8	35.4 $\pm$ 0.2	25.1 $\pm$ 0.8	22.6 $\pm$ 0.2	28.4 $\pm$ 0.7	19.4 $\pm$ 0.2	26.8 $\pm$ 0.3	24.4 $\pm$ 0.1	25.3 $\pm$ 0.1	26.0 $\pm$ 0.8	27.2 $\pm$ 0.4
	SAS	40.8 $\pm$ 0.3	26.6 $\pm$ 0.9	32.7 $\pm$ 0.2	19.9 $\pm$ 0.7	<u>32.5 <math>\pm</math>0.7</u>	<u>33.1 <math>\pm</math>0.1</u>	<u>28.7 <math>\pm</math>0.4</u>	<u>30.6 <math>\pm</math>0.3</u>	28.8 $\pm$ 0.2	26.7 $\pm$ 0.1	29.3 $\pm$ 0.2	28.0 $\pm$ 0.4
	DRA	<u>45.7 <math>\pm</math>0.2</u>	<u>34.8 <math>\pm</math>0.9</u>	<u>40.2 <math>\pm</math>0.2</u>	<u>28.4 <math>\pm</math>0.7</u>	29.3 $\pm$ 0.7	31.4 $\pm$ 0.5	21.7 $\pm$ 0.2	26.0 $\pm$ 0.2	21.7 $\pm$ 0.1	22.4 $\pm$ 0.0	23.2 $\pm$ 0.9	24.2 $\pm$ 0.9
	DTR (Ours)	<b>49.0 <math>\pm</math>0.1</b>	<b>37.4 <math>\pm</math>0.2</b>	<b>42.7 <math>\pm</math>0.3</b>	<b>30.4 <math>\pm</math>0.5</b>	<b>35.9 <math>\pm</math>0.0</b>	<b>36.7 <math>\pm</math>0.1</b>	<b>30.8 <math>\pm</math>0.3</b>	<b>32.7 <math>\pm</math>0.1</b>	<b>1.3 <math>\pm</math>0.1</b>	<b>0.9 <math>\pm</math>0.4</b>	<b>2.4 <math>\pm</math>0.0</b>	<b>2.2 <math>\pm</math>0.2</b>

Table 2: Main results on response quality and latency for two STMD-Sim benchmarks across three LLM backbones. We report first-token (FTL) and end-to-end (E2EL) latency in seconds (Sec.). Results are reported as "mean  $\pm$  standard deviation" over three random seeds. Best results are in **bold** and second-best results are underlined.

with Search) (Press et al., 2023), which decomposes query into a small set of follow-up questions, performs web search for each sub-question, and synthesizes a grounded answer from the aggregated evidence; and (iv) DRA (Deep Research Agent) (Yamada et al., 2025a), which follows an open, reproducible research-agent workflow with planning, iterative retrieval, evidence aggregation, and long-form synthesis.

### 5.3 Tasks and Metrics

For automatic evaluation, we measure response quality and latency following common practice (Rudd et al., 2025; NVIDIA, 2025). For each project, we evaluate on the held-out online meeting and use the first 10 meetings for warmup and memory construction. We sample key turns (topic transitions and conclusion-forming turns), form a query from the online context, and summarize a reference answer from the project’s golden conclusion. We report mean and standard deviation over three random seeds. **Response Quality** is measured by ROUGE-1, ROUGE-L, BLEU-1, and BLEU-N, and **Response Latency** by first-token and end-to-end latency in seconds. Following Wang et al. (2023b), we also run human-involved studies with three expert biomedical teams covering cancer, sepsis, and cardiometabolic disease, with details in Section 5.5 and Appendix B.2.

### 5.4 Main Results

We evaluate DTR across three LLM backbones and representative baselines, measuring both re-

sponse quality and response latency (Table 2). DTR improves response quality while maintaining low latency, whereas all baselines require substantially higher latency to reach comparable quality in some settings. For example, on gpt-oss-120b, DTR achieves the best response quality on both STMD<sub>1</sub> and STMD<sub>2</sub> (e.g., ROUGE-1 of 52.1 on STMD<sub>1</sub> and 26.8 on STMD<sub>2</sub>). It also reduces latency to 1.0/1.1 seconds FTL and 1.9/2.1 seconds E2EL, compared with 9.7/8.1 seconds FTL and 11.1/9.5 seconds E2EL under Standard method, and 19.6–24.7 seconds FTL for DRA, ReAct, and SAS. This quality–latency advantage remains stable on Gemma-3-12B-it and Llama-3.1-8B-Instruct. On these backbones, DTR achieves leading ROUGE scores and the best BLEU-based scores, while keeping FTL within 0.9–1.3 seconds and E2EL within 1.9–2.4 seconds across STMD<sub>1</sub> and STMD<sub>2</sub>.

### 5.5 Human Involved Validation

In addition, we conduct human-involved usability testing<sup>5</sup> to assess response quality under real Team-AI interaction. The human evaluation includes three biomedical research teams covering cancer, sepsis, and cardio-metabolic disease. In each meeting session, participants provide point-wise ratings on a 1-5 scale along three dimensions, including response *Quality*, *Helpfulness* to the project goal, and *Personalization* to the recipient team member, using the rubric in Table 7. As

<sup>5</sup>Details of the participant information and evaluation protocol are provided in Appendix B.2.

LLM Backbone	Method	Response Quality (%) $\uparrow$								Response Latency (Sec.) $\downarrow$			
		ROUGE-1		ROUGE-L		BLEU-1		BLEU-N		FTL		E2EL	
		STMD <sub>1</sub>	STMD <sub>2</sub>	STMD <sub>1</sub>	STMD <sub>2</sub>	STMD <sub>1</sub>	STMD <sub>2</sub>	STMD <sub>1</sub>	STMD <sub>2</sub>	STMD <sub>1</sub>	STMD <sub>2</sub>	STMD <sub>1</sub>	STMD <sub>2</sub>
Llama-3.1-8B-Instruct	Standard	38.8 $\pm$ 0.2	29.8 $\pm$ 0.6	31.6 $\pm$ 0.2	25.5 $\pm$ 0.7	10.8 $\pm$ 0.8	20.0 $\pm$ 0.7	9.3 $\pm$ 0.7	18.5 $\pm$ 0.3	7.6 $\pm$ 0.2	8.1 $\pm$ 0.0	8.2 $\pm$ 0.2	8.9 $\pm$ 0.2
	w/o Injection	42.1 $\pm$ 0.3	33.5 $\pm$ 0.1	35.4 $\pm$ 0.4	26.7 $\pm$ 0.2	32.1 $\pm$ 0.7	32.8 $\pm$ 0.4	28.9 $\pm$ 0.1	29.1 $\pm$ 0.2	1.2 $\pm$ 0.2	1.4 $\pm$ 0.4	2.1 $\pm$ 0.3	2.6 $\pm$ 0.2
	w/o ConFilter	45.2 $\pm$ 0.5	35.1 $\pm$ 0.3	38.9 $\pm$ 0.3	27.8 $\pm$ 0.2	28.5 $\pm$ 0.1	29.4 $\pm$ 0.3	25.6 $\pm$ 0.6	26.2 $\pm$ 0.2	1.8 $\pm$ 0.1	1.5 $\pm$ 0.2	3.1 $\pm$ 0.3	2.9 $\pm$ 0.1
	w/o GateLLM	46.8 $\pm$ 0.1	36.2 $\pm$ 0.2	40.1 $\pm$ 0.2	28.9 $\pm$ 0.3	22.3 $\pm$ 0.4	25.6 $\pm$ 0.5	19.8 $\pm$ 0.2	22.7 $\pm$ 0.4	2.2 $\pm$ 0.1	1.8 $\pm$ 0.5	3.6 $\pm$ 0.2	3.1 $\pm$ 0.3
	DTR (Ours)	49.0 $\pm$ 0.1	37.4 $\pm$ 0.2	42.7 $\pm$ 0.3	30.4 $\pm$ 0.5	35.9 $\pm$ 0.0	36.7 $\pm$ 0.1	30.8 $\pm$ 0.3	32.7 $\pm$ 0.1	1.3 $\pm$ 0.1	0.9 $\pm$ 0.4	2.4 $\pm$ 0.0	2.2 $\pm$ 0.2

Table 3: Comparison with variants of DTR on Llama-3.1-8B-Instruct. w/o Injection removes memory and persona injection during offline meeting simulation. w/o ConFilter disables confidence-based filtering. w/o GateLLM disables the gate LLM and uses a random meeting round for offline search.

Method	Quality $\uparrow$	Helpfulness $\uparrow$	Personalization $\uparrow$	FTL $\downarrow$
GPT-5.2-Thinking	3.25	3.19	3.13	35.5
DTR (Ours)	4.03	4.02	4.23	1.99

Table 4: Human evaluation results comparing DTR and GPT-5.2-Thinking on Quality, Helpfulness, Personalization, and first-token latency (FTL) in seconds.

shown in Table 4, averaged over 200 feedback instances, DTR (Qwen-Plus API backbone (Yang et al., 2025)) outperforms GPT-5.2-Thinking on all dimensions, achieving higher Quality (4.03 vs. 3.25), Helpfulness (4.02 vs. 3.19), and Personalization (4.23 vs. 3.13), while also reducing first-token latency from 35.5s to 1.99s.

## 5.6 Ablation Study

### 5.6.1 Variants Comparison

We evaluate DTR with ablations to isolate each component’s contribution (Table 3). DTR performs best on both STMD-Sim benchmarks, outperforming Standard and all ablations on response-quality metrics. Among variants, w/o GateLLM achieves the highest ROUGE but still trails DTR on STMD<sub>1</sub>/STMD<sub>2</sub> (e.g., 46.8/36.2 vs. 49.0/37.4 ROUGE-1), suggesting that gating improves when to search and what to distill during rehearsal. Removing memory/persona injection reduces ROUGE-1 to 42.1/33.5. Disabling confidence-based filtering also lowers ROUGE on both benchmarks, consistent with filtering out low-confidence rehearsal meetings and improving distilled-knowledge reliability. Overall, each component contributes to DTR’s gains.

### 5.6.2 Effect of the Number of Offline Rehearsal Meetings

We analyze the effect of the offline rehearsal budget on the Llama-3.1-8B-Instruct backbone by varying the number of offline rehearsal meetings used to construct the rehearsal knowledge database. Fig-

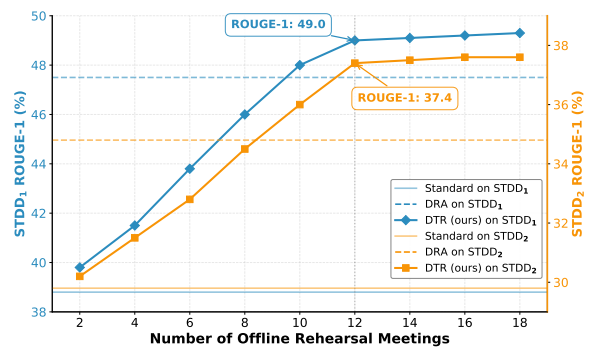


Figure 4: Effect of the number of offline rehearsal meetings on ROUGE-1 for DTR and baselines on STMD<sub>1</sub> and STMD<sub>2</sub> with Llama-3.1-8B-Instruct.

ure 4 reports ROUGE-1 on STMD<sub>1</sub> and STMD<sub>2</sub>. Increasing the number of rehearsal meetings consistently improves DTR on both benchmarks, with gains that taper at larger budgets and saturate at 49.0 ROUGE-1 on STMD<sub>1</sub> and 37.4 ROUGE-1 on STMD<sub>2</sub>. DTR benefits from additional rehearsal evidence, with diminishing returns beyond a moderate rehearsal budget.

## 6 Conclusion and Future Work

We propose DTR, a rehearsal-based framework that enables deep thinking in scientific team meetings without prohibitive online latency. DTR decouples multi-step reasoning and knowledge collection from synchronous interaction by shifting to an offline rehearsal phase. In addition, we introduce STMD, a hybrid benchmark comprising authentic transcripts from three biomedical teams and simulated deliberations synthesized from PubMed literature. Empirical experiments demonstrate DTR improves response quality while reducing inference latency compared with baselines. Future work will extend DTR to other latency-sensitive domains and explore adaptive rehearsal strategies for personalized, multimodal Team-AI collaboration.

## 561 Limitations

562 While DTR demonstrates promising results in en-  
563 abling low-latency, high-quality human–AI inter-  
564 actions, we acknowledge several limitations of the  
565 current work:

566 • **Knowledge Coverage Limits.** Because rehearsal  
567 simulates anticipated meeting scenarios in advance,  
568 it may not provide complete coverage of all dis-  
569 cussion directions that arise at deployment time.  
570 Simulated meetings may miss certain questions,  
571 omit relevant evidence, or express information  
572 in ways that reduce later retrievability. Future  
573 work can strengthen coverage while preserving  
574 low latency by improving simulation diversity,  
575 adding lightweight verification, and incorporating  
576 uncertainty-aware fallback behaviors.

577 • **Simplified Tool-use Assumptions.** Our study  
578 focuses on two widely used capabilities in cur-  
579 rent LLM assistants: multi-step reasoning and web  
580 search. This choice keeps the rehearsal setting  
581 controlled and reproducible, but it simplifies the  
582 broader tool landscape that may arise in real work-  
583 flows, such as querying specialized biomedical  
584 databases, calling domain-specific analysis tools,  
585 or interacting with lab software. In addition, we do  
586 not incorporate multimodal inputs such as figures,  
587 tables, slides, or imaging artifacts, which are often  
588 central to scientific collaboration. Extending re-  
589 hearsal to broader tool ecosystems and multimodal  
590 evidence is a promising direction.

591 • **Offline Compute Budget.** Our method introduces  
592 additional offline computation for rehearsal simula-  
593 tion, including tool-assisted retrieval and reasoning  
594 during meeting generation. This upfront cost is  
595 deliberately incurred to eliminate latency and fric-  
596 tion in the human–AI interaction loop during live  
597 deployment, which is the primary design objective.  
598 In practice, the offline budget is adjustable. We  
599 can vary the number of simulated rehearsals per  
600 project to balance preparation cost against down-  
601 stream performance, allowing the approach flexible  
602 across different computational constraints.

603 • **Inherent Latency in Real-world Interfaces.** Al-  
604 though DTR reduces online latency compared with  
605 baselines that perform multi-step reasoning and  
606 web search during inference, deployment in speech-  
607 based meetings typically requires automatic speech  
608 recognition (ASR) for transcription and text-to-  
609 speech (TTS) for speech output, which can add  
610 end-to-end system latency. Future work will ex-  
611 plore streaming ASR and TTS integration and in-

cremental processing to reduce interface overhead  
while preserving response quality.

## Ethics Statement

Real-world data collection and human evaluation  
in this study involve three biomedical research  
teams, each comprising one Principal Investigator  
(PI) and five PhD student researchers. The teams  
are recruited from established biomedical research  
institutions and cover cancer, sepsis, and cardio-  
metabolic disease. We also recruit two PhD stu-  
dents with biomedical backgrounds to conduct the  
data quality assessment. Participants are compen-  
sated at their standard institutional research rates  
for time spent in the experimental sessions. To  
reduce potential bias, participants are blinded to  
the study hypotheses, model architecture, and ex-  
perimental conditions throughout the evaluation.  
Informed consent is obtained prior to participa-  
tion, and participants agree that their interaction  
data may be anonymized and used for research  
purposes.

The experimental sessions use simulated sci-  
entific meetings grounded in publicly available  
PubMed literature. The study does not involve  
personal, sensitive, or confidential information. All  
procedures follow institutional guidelines, and no  
additional ethics board review is required under  
the ACL Ethics Policy. Upon publication, we will  
release the dataset, experimental protocols, and an-  
notation guidelines to support reproducibility and  
future research on AI-assisted scientific collabora-  
tion systems.

## References

- Zhenni Bi, Kai Han, Chuanjian Liu, Yehui Tang, and  
Yunhe Wang. 2024. Forest-of-thought: Scaling test-  
time compute for enhancing llm reasoning. *arXiv  
preprint arXiv:2412.09078*.
- Ruizhi Cheng, Surendra Pathak, Guowu Xie, Matteo  
Varvello, Songqing Chen, and Bo Han. 2025. Hello,  
genai? dissecting human to generative ai calling. In  
*Proceedings of the 2025 ACM Internet Measurement  
Conference*, pages 308–324.
- cyanheads. 2025. [pubmed-mcp-server: An mcp server  
for pubmed via ncbi e-utilities](#). GitHub repository.  
Accessed 2026-01-01.
- Yichao Fu, Xuewei Wang, Yuandong Tian, and Jiawei  
Zhao. 2025. Deep think with confidence. *arXiv  
preprint arXiv:2508.15260*.

660	Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon,	Romain Lacombe, Kerrie Wu, and Eddie Dilworth.	714
661	Pengfei Liu, Yiming Yang, Jamie Callan, and Gra-	2025. Don't think twice! over-reasoning im-	715
662	ham Neubig. 2023. Pal: Program-aided language	pairs confidence calibration. <i>arXiv preprint</i>	716
663	models. In <i>International Conference on Machine</i>	<i>arXiv:2508.15050</i> .	717
664	<i>Learning</i> , pages 10764–10799. PMLR.		
665	Silin Gao, Jane Dwivedi-Yu, Ping Yu, Xiaoqing Ellen	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	718
666	Tan, Ramakanth Pasunuru, Olga Golovneva, Kous-	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	719
667	tuv Sinha, Asli Celikyilmaz, Antoine Bosselut, and	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	720
668	Tianlu Wang. 2025. Efficient tool use with chain-	täschel, et al. 2020. Retrieval-augmented generation	721
669	of-abstraction reasoning. In <i>Proceedings of the 31st</i>	for knowledge-intensive nlp tasks. <i>Advances in neu-</i>	722
670	<i>International Conference on Computational Linguis-</i>	<i>ral information processing systems</i> , 33:9459–9474.	723
671	<i>tics</i> , pages 2727–2743.		
672	Google. 2025. <a href="#">google/gemma-3-12b-it: Model card.</a>	Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii	724
673	Hugging Face model card. Accessed 2026-01-01.	Khizbullin, and Bernard Ghanem. 2023. Camel:	725
674	Google DeepMind. 2025. <a href="#">Gemini 3 pro: Model card.</a>	Communicative agents for "mind" exploration of	726
675	Model card, Google DeepMind. Accessed 2025-12-	large language model society. <i>Advances in Neural</i>	727
676	27.	<i>Information Processing Systems</i> , 36:51991–52008.	728
677	Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and	Lizi Liao, Grace Hui Yang, and Chirag Shah. 2023.	729
678	Timothy Lillicrap. 2023. Mastering diverse do-	Proactive conversational agents. In <i>Proceedings of</i>	730
679	domains through world models. <i>arXiv preprint</i>	<i>the Sixteenth ACM International Conference on Web</i>	731
680	<i>arXiv:2301.04104</i> .	<i>Search and Data Mining</i> , pages 1244–1247.	732
681	Tao He, Lizi Liao, Yixin Cao, Yuanxing Liu, Ming	Jung Hoon Lim, Sunjae Kwon, Zonghai Yao, John P	733
682	Liu, Zerui Chen, and Bing Qin. 2024. Planning like	Lalor, and Hong Yu. 2024. Large language model-	734
683	human: A dual-process framework for dialogue plan-	based role-playing for personalized medical jargon	735
684	ning. <i>arXiv preprint arXiv:2406.05374</i> .	extraction. <i>arXiv preprint arXiv:2408.05555</i> .	736
685	Tao He, Lizi Liao, Ming Liu, and Bing Qin. 2025. Sim-	Xinyan Liu, Jiaxin Lu, and Ang Xu. 2025. Dynamic	737
686	ulating before planning: Constructing intrinsic user	agenda-aware real-time meeting summarization with	738
687	world model for user-tailored dialogue policy plan-	large language models. <i>Journal of King Saud Univer-</i>	739
688	ning. In <i>Proceedings of the 48th International ACM</i>	<i>sity Computer and Information Sciences</i> , 37(9):1–22.	740
689	<i>SIGIR Conference on Research and Development in</i>	Yaxi Lu, Shenzhi Yang, Cheng Qian, Guirong Chen,	741
690	<i>Information Retrieval</i> , pages 645–655.	Qinyu Luo, Yesai Wu, Huadong Wang, Xin Cong,	742
691	Ehud Karpas, Omri Abend, Yonatan Belinkov, Barak	Zhong Zhang, Yankai Lin, et al. 2024. Proactive	743
692	Lenz, Opher Lieber, Nir Ratner, Yoav Shoham, Hofit	agent: Shifting llm agents from reactive responses to	744
693	Bata, Yoav Levine, Kevin Leyton-Brown, et al. 2022.	active assistance. <i>arXiv preprint arXiv:2410.12361</i> .	745
694	Mrkl systems: A modular, neuro-symbolic architec-	Mykola Maslych, Mohammadreza Katebi, Christo-	746
695	ture that combines large language models, external	pher Lee, Yahya Hmaiti, Amirpouya Ghasemaghaei,	747
696	knowledge sources and discrete reasoning. <i>arXiv</i>	Christian Pumarada, Janneese Palmer, Esteban	748
697	<i>preprint arXiv:2205.00445</i> .	Segarra Martinez, Marco Emporio, Warren Snipes,	749
698	Namyoun Kim, Kai Tzu-iunn Ong, Yeonjun Hwang,	et al. 2025. Mitigating response delays in free-form	750
699	Minseok Kang, Iiseo Jihn, Gayoung Kim, Minju Kim,	conversations with llm-powered intelligent virtual	751
700	and Jinyoung Yeo. 2025. Principles: Synthetic strat-	agents. In <i>Proceedings of the 7th ACM Conference</i>	752
701	egy memory for proactive dialogue agents. <i>arXiv</i>	<i>on Conversational User Interfaces</i> , pages 1–15.	753
702	<i>preprint arXiv:2509.17459</i> .	Meta. 2024. <a href="#">meta-llama/llama-3.1-8b-instruct: Model</a>	754
703	Pavel Král, Věra Králová, and Petr Šimáček. 2023. The	<a href="#">card</a> . Hugging Face model card. Accessed 2026-01-	755
704	impact of interactions before, during and after meet-	01.	756
705	ings on meeting effectiveness: a coordination the-	NVIDIA. 2025. <a href="#">Llm inference benchmarking: Funda-</a>	757
706	ory perspective. <i>Measuring Business Excellence</i> ,	<a href="#">mental concepts</a> . NVIDIA Technical Blog. Accessed	758
707	27(3):403–420.	2026-01-01.	759
708	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying	OpenAI. 2024. <a href="#">GPT-4o System Card</a> . Technical report,	760
709	Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.	OpenAI. Accessed: 2025-12-29.	761
710	Gonzalez, Hao Zhang, and Ion Stoica. 2023. Ef-	OpenAI. 2025a. <a href="#">GPT-5 System Card</a> . Technical report,	762
711	efficient memory management for large language	OpenAI. Accessed: 2025-12-29.	763
712	model serving with pagedattention. <i>arXiv preprint</i>	OpenAI. 2025b. <a href="#">gpt-oss-120b &amp; gpt-oss-20b model</a>	764
713	<i>arXiv:2309.06180</i> .	<a href="#">card</a> . OpenAI model card. Accessed 2026-01-01.	765

766	OpenAI. 2025c. <a href="#">Update to gpt-5 system card: Gpt-5.2</a> .	822
767	System card addendum, OpenAI. Accessed 2025-12-	823
768	27.	824
769	Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith	825
770	Ringel Morris, Percy Liang, and Michael S Bernstein. 2023.	826
771	Generative agents: Interactive simulacra	827
772	of human behavior. In <i>Proceedings of the 36th annual</i>	828
773	<i>acm symposium on user interface software and</i>	829
774	<i>technology</i> , pages 1–22.	830
775	Adam Paszke, Sam Gross, Francisco Massa, Adam	831
776	Lerer, James Bradbury, Gregory Chanan, Trevor	832
777	Killeen, Zeming Lin, Natalia Gimelshein, Luca	833
778	Antiga, et al. 2017. Automatic differentiation in py-	834
779	torch. In <i>Advances in Neural Information Processing</i>	835
780	<i>Systems</i> .	836
781	Hannah Pelikan and Emily Hofstetter. 2023. Managing	837
782	delays in human-robot interaction. <i>ACM Transac-</i>	838
783	<i>tions on Computer-Human Interaction</i> , 30(4):1–42.	839
784	Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt,	840
785	Noah A Smith, and Mike Lewis. 2023. Measuring	841
786	and narrowing the compositionality gap in language	842
787	models. In <i>Findings of the Association for Computa-</i>	843
788	<i>tional Linguistics: EMNLP 2023</i> , pages 5687–5711.	844
789	Michael H Prince, Henry Chan, Aikaterini Vriza, Tao	845
790	Zhou, Varuni K Sastry, Yanqi Luo, Matthew T	846
791	Dearing, Ross J Harder, Rama K Vasudevan, and	847
792	Mathew J Cherukara. 2024. Opportunities for re-	848
793	trieval and tool augmented large language models	849
794	in scientific facilities. <i>npj Computational Materials</i> ,	850
795	10(1):251.	851
796	Biqing Qi, Kaiyan Zhang, Kai Tian, Haoxiang	852
797	Li, Zhang-Ren Chen, Sihang Zeng, Ermo Hua,	853
798	Hu Jinfang, and Bowen Zhou. 2024. Large lan-	854
799	guage models as biomedical hypothesis genera-	855
800	tors: a comprehensive evaluation. <i>arXiv preprint</i>	856
801	<i>arXiv:2407.08940</i> .	857
802	Ethan M. Rudd, Christopher Andrews, and Philip Tully.	858
803	2025. A practical guide for evaluating llms and llm-	859
804	reliant systems. <i>arXiv preprint arXiv:2506.13023</i> .	860
805	Vishal Vivek Saley, Goonjan Saha, Rocktim Jyoti	861
806	Das, Dinesh Raghu, et al. 2024. Meditod: An en-	862
807	glish dialogue dataset for medical history taking	863
808	with comprehensive annotations. <i>arXiv preprint</i>	864
809	<i>arXiv:2410.14204</i> .	865
810	Eric W Sayers, Jeff Beck, Evan E Bolton, J Rodney	866
811	Brister, Jessica Chan, Donald C Comeau, Ryan Con-	867
812	nor, Michael DiCuccio, Catherine M Farrell, Michael	868
813	Feldgarden, et al. 2024. <a href="#">Database resources of the national</a>	869
814	<a href="#">center for biotechnology information</a> . <i>Nucleic</i>	870
815	<i>Acids Research</i> , 52(D1):D33–D43.	871
816	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta	872
817	Raileanu, Maria Lomeli, Eric Hambro, Luke Zettle-	873
818	moyer, Nicola Cancedda, and Thomas Scialom. 2023.	874
819	Toolformer: Language models can teach themselves	875
820	to use tools. <i>Advances in Neural Information Pro-</i>	
821	<i>cessing Systems</i> , 36:68539–68551.	
	Jan B Schmutz, Neal Outland, Sophie Kerstan, Eleni	
	Georganta, and Anna-Sophie Ulfert. 2024. Ai-	
	teaming: Redefining collaboration in the digital era.	
	<i>Current Opinion in Psychology</i> , 58:101837.	
	Julian Schrittwieser, Ioannis Antonoglou, Thomas Hu-	
	bert, Karen Simonyan, Laurent Sifre, Simon Schmitt,	
	Arthur Guez, Edward Lockhart, Demis Hassabis,	
	Thore Graepel, et al. 2020. Mastering atari, go, chess	
	and shogi by planning with a learned model. <i>Nature</i> ,	
	588(7839):604–609.	
	Amanpreet Singh, Joseph Chee Chang, Chloe Anastasi-	
	ades, Dany Haddad, Aakanksha Naik, Amber Tanaka,	
	Angele Zamarron, Cecile Nguyen, Jena D Hwang,	
	Jason Dunkleberger, et al. 2025. Ai2 scholar qa: Or-	
	ganized literature synthesis with attribution. <i>arXiv</i>	
	<i>preprint arXiv:2504.10861</i> .	
	Nishant Subramani, Jason Eisner, Justin Svegliato, Ben-	
	jamin Van Durme, Yu Su, and Sam Thomson. 2025.	
	Mice for cats: Model-internal confidence estima-	
	tion for calibrating agents with tools. <i>arXiv preprint</i>	
	<i>arXiv:2504.20168</i> .	
	Yufei Tao, Ameeta Agrawal, Judit Dombi, Tetyana	
	Sydorenko, and Jung In Lee. 2024. Chatgpt role-	
	play dataset: Analysis of user motives and model	
	naturalness. <i>arXiv preprint arXiv:2403.18121</i> .	
	Amir Taubenfeld, Tom Sheffer, Eran Ofek, Amir Feder,	
	Ariel Goldstein, Zorik Gekhman, and Gal Yona. 2025.	
	Confidence improves self-consistency in llms. <i>arXiv</i>	
	<i>preprint arXiv:2502.06233</i> .	
	Guangzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Man-	
	dlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and An-	
	ima Anandkumar. 2023a. Voyager: An open-ended	
	embodied agent with large language models. <i>arXiv</i>	
	<i>preprint arXiv:2305.16291</i> .	
	Hongru Wang, Rui Wang, Fei Mi, Zezhong Wang,	
	Ruifeng Xu, and Kam-Fai Wong. 2023b. Chain-	
	of-thought prompting for responding to in-depth	
	dialogue questions with llm. <i>arXiv preprint</i>	
	<i>arXiv:2305.11792</i> .	
	Huaijie Wang, Shibo Hao, Hanze Dong, Shenao Zhang,	
	Yilin Bao, Ziran Yang, and Yi Wu. 2025. Offline	
	reinforcement learning for llm multi-step reasoning.	
	In <i>Findings of the Association for Computational</i>	
	<i>Linguistics: ACL 2025</i> , pages 8881–8893.	
	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le,	
	Ed Chi, Sharan Narang, Aakanksha Chowdhery, and	
	Denny Zhou. 2022. Self-consistency improves chain	
	of thought reasoning in language models. <i>arXiv</i>	
	<i>preprint arXiv:2203.11171</i> .	
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	
	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	
	et al. 2022. Chain-of-thought prompting elicits rea-	
	soning in large language models. <i>Advances in neural</i>	
	<i>information processing systems</i> , 35:24824–24837.	

876	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	from chain-of-thought reasoning to language agents.	932
877	Chaumond, Clement Delangue, Anthony Moi, Pier-	<i>ACM Computing Surveys</i> , 57(8):1–39.	933
878	ric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz,		
879	et al. 2020. Transformers: State-of-the-art natural	Tianshi Zheng, Zheyue Deng, Hong Ting Tsang, Weiqi	934
880	language processing. In <i>Proceedings of the 2020 con-</i>	Wang, Jiabin Bai, Zihao Wang, and Yangqiu Song.	935
881	<i>ference on empirical methods in natural language</i>	2025. From automation to autonomy: A survey on	936
882	<i>processing: system demonstrations</i> , pages 38–45.	large language models in scientific discovery. <i>arXiv</i>	937
		<i>preprint arXiv:2505.13259</i> .	938
883	Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas		
884	Muennighoff. 2023. <b>C-pack: Packaged resources</b>	Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman,	939
885	<b>to advance general chinese embedding</b> .	Haohan Wang, and Yu-Xiong Wang. 2023. Lan-	940
		guage agent tree search unifies reasoning acting	941
886	Ikuya Yamada, Wataru Ikeda, Ko Yoshida, Mengyu Ye,	and planning in language models. <i>arXiv preprint</i>	942
887	Hinata Sugimoto, Masatoshi Suzuki, Hisanori Ozaki,	<i>arXiv:2310.04406</i> .	943
888	and Jun Suzuki. 2025a. An open and reproducible		
889	deep research agent for long-form question answer-		
890	ing. <i>arXiv preprint arXiv:2512.13059</i> .		
891	Yutaro Yamada, Robert Tjarko Lange, Cong Lu, Shen-		
892	gran Hu, Chris Lu, Jakob Foerster, Jeff Clune, and		
893	David Ha. 2025b. The ai scientist-v2: Workshop-		
894	level automated scientific discovery via agentic tree		
895	search. <i>arXiv preprint arXiv:2504.08066</i> .		
896	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,		
897	Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,		
898	Chengen Huang, Chenxu Lv, et al. 2025. Qwen3		
899	technical report. <i>arXiv preprint arXiv:2505.09388</i> .		
900	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,		
901	Tom Griffiths, Yuan Cao, and Karthik Narasimhan.		
902	2023. Tree of thoughts: Deliberate problem solving		
903	with large language models. <i>Advances in neural</i>		
904	<i>information processing systems</i> , 36:11809–11822.		
905	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak		
906	Shafran, Karthik R Narasimhan, and Yuan Cao. 2022.		
907	React: Synergizing reasoning and acting in language		
908	models. In <i>The eleventh international conference on</i>		
909	<i>learning representations</i> .		
910	Guangtao Zeng, Wenmian Yang, Zeqian Ju, Yue Yang,		
911	Sicheng Wang, Ruisi Zhang, Meng Zhou, Jiaqi Zeng,		
912	Xiangyu Dong, Ruoyu Zhang, et al. 2020. Med-		
913	dialog: Large-scale medical dialogue datasets. In		
914	<i>Proceedings of the 2020 conference on empirical</i>		
915	<i>methods in natural language processing (EMNLP)</i> ,		
916	pages 9241–9250.		
917	Ceyao Zhang, Kaijie Yang, Siyi Hu, Zihao Wang,		
918	Guanghe Li, Yihang Sun, Cheng Zhang, Zhaowei		
919	Zhang, Anji Liu, Song-Chun Zhu, et al. 2024. Proa-		
920	gent: building proactive cooperative agents with large		
921	language models. In <i>Proceedings of the AAAI Con-</i>		
922	<i>ference on Artificial Intelligence</i> , volume 38, pages		
923	17591–17599.		
924	Kai Zhang, Yangyang Kang, Fubang Zhao, and Xi-		
925	aozhong Liu. 2023. Llm-based medical assistant		
926	personalization with short-and long-term memory		
927	coordination. <i>arXiv preprint arXiv:2309.11696</i> .		
928	Zhuosheng Zhang, Yao Yao, Aston Zhang, Xiangru		
929	Tang, Xinbei Ma, Zhiwei He, Yiming Wang, Mark		
930	Gerstein, Rui Wang, Gongshen Liu, et al. 2025. Ig-		
931	nitening language intelligence: The hitchhiker’s guide		

## A Data Statistics

We select PubMed papers (Sayers et al., 2024) spanning cancer, Alzheimer’s disease, and sepsis, prioritizing papers published between January 1, 2024 and January 1, 2025 with concrete methodological descriptions and sufficient technical detail for project specification. The raw corpus includes 443 cancer papers, 198 Alzheimer’s papers, and 41 sepsis papers. As described in Section 4, we treat each paper as one project and derive three paper-grounded artifacts: a compact mini-proposal, a golden conclusion, and an ordered topic sequence that guides meeting progression. Using the same construction pipeline, we build two backbone-specific variants for STMD-Sim: STMD<sub>1</sub> (GPT-5 (OpenAI, 2025a)) and STMD<sub>2</sub> (GPT-4o (OpenAI, 2024)). STMD<sub>1</sub> contains 300 teams, producing 3,300 meetings and 39,600 utterances, with an average of 74 tokens per utterance and 811 evaluation queries. STMD<sub>2</sub> contains 2,046 teams, producing 22,506 meetings and 270,072 utterances, with an average of 93 tokens per utterance and 4,297 evaluation queries. In both variants, each team has 10 historical meetings and one held-out online meeting, each with 12 utterances. In addition, STMD-Real contains 12 weeks of multi-party lab meeting transcripts from three biomedical research teams, each comprising one PI and five PhD researchers. The first 10 weeks yield 10 historical meetings per team, and the remaining 2 weeks yield 2 online evaluation meetings per team. The detailed statistics are provided in Table 5.

## B Human Evaluation Details

### B.1 Dataset Quality Assessment

**Human Experts** We recruit two PhD researchers with biomedical backgrounds to assess the quality of the constructed STMD-Sim meetings. They are compensated at their standard institutional research rates for time spent in the annotation sessions. Informed consent is obtained prior to participation, and they agree that their annotation data may be anonymized and used for research purposes.

**Assessment Metrics** We randomly sample 30 meetings from each dataset variant for independent expert annotation. Using the rubric in Table 6, experts rate each meeting on a 1–5 scale along three criteria, including *Meeting Coherence*, *Factual Faithfulness*, and *Role Consistency*. We report average quality scores and inter-annotator agreement mea-

<i>Raw Data Statistics of PubMed Papers</i>		
# Cancer Papers	443	
# Alzheimer’s Papers	198	
# Sepsis Papers	41	
<i>Synthetic Data Statistics (STMD-Sim)</i>		
	STMD <sub>1</sub>	STMD <sub>2</sub>
# Teams	300	2046
# Historical Meetings / Team	10	10
# Online Meetings / Team	1	1
# Utterances / Meeting	12	12
# Avg. Tokens / Utterance	74	93
# Total Utterances	39,600	270,072
# Total Meetings	3,300	22,506
# Total Evaluation Queries	811	4,297
<i>Real-World Data Statistics (STMD-Real)</i>		
Collection Period	12 Weeks	
# Teams	3	
# Members / Team	6	
# Roles / Team	1 PI + 5 PhD	
# Historical Meetings / Team	10	
# Online Meetings / Team	2	
# Avg. Utterances / Meeting	47	
# Total Meetings	36	
# Total Utterances	1,682	

Table 5: Statistics of the constructed STMD dataset.

sured by Cohen’s  $\kappa$  in Table 1.

### B.2 Human-Involved Model Evaluation

**Human Participants** We recruit three biomedical research teams for human-involved evaluation, each comprising one PI and five PhD student researchers. The teams cover cancer, sepsis, and cardio-metabolic disease. Participants are compensated at their standard institutional research rates for time spent in the sessions. To reduce potential bias, participants are blinded to the study hypotheses and experimental conditions throughout the evaluation. Informed consent is obtained prior to participation, and participants agree that their interaction data may be anonymized and used for research purposes.

**Evaluation Metrics** In each meeting session, participants provide point-wise ratings on a 1–5 scale along three dimensions, including response *Quality*, *Helpfulness* to the project goal, and *Personalization* to the recipient team member, following the rubric in Table 7. We report the average ratings over all feedback instances in Table 4.

<b>Meeting Coherence</b>		
Score	Label	Description
1	Poor Coherence	Meeting is disjointed, with abrupt topic shifts and broken turn-to-turn logic.
2	Weak Coherence	Partly connected, but frequent gaps, unclear references, or uneven topic progression.
3	Acceptable Coherence	Generally coherent, though some jumps or weak connections reduce smoothness.
4	Good Coherence	Coherent and well-structured, with clear topic flow and sensible turn transitions.
5	Excellent Coherence	Highly coherent, with smooth progression and strong logical connections throughout.
<b>Factual Faithfulness</b>		
Score	Label	Description
1	Poor Faithfulness	Many claims conflict with the source paper or are clearly unsupported.
2	Weak Faithfulness	Some statements are plausible, but key details are incorrect or insufficiently supported.
3	Acceptable Faithfulness	Mostly faithful, with minor inaccuracies or vague references to paper evidence.
4	Good Faithfulness	Faithful to the paper, with only minor missing details or imprecise wording.
5	Excellent Faithfulness	Highly faithful, accurately reflecting paper evidence, results, and conclusions.
<b>Role Consistency</b>		
Score	Label	Description
1	Poor Consistency	Roles are confused or interchangeable, with frequent role drift across turns.
2	Weak Consistency	Some role alignment, but noticeable drift in expertise, tone, or responsibilities.
3	Acceptable Consistency	Roles are mostly consistent, with occasional drift or mismatched level of detail.
4	Good Consistency	Roles are consistent and distinguishable, matching expected expertise and focus.
5	Excellent Consistency	Roles are strongly consistent, with clear specialization and stable role behavior.

Table 6: Data quality scoring form.

<b>Quality</b>		
Score	Label	Description
1	Poor Quality	Incorrect or incoherent, largely irrelevant, with major factual issues.
2	Weak Quality	Partially correct but vague, with omissions that reduce reliability.
3	Acceptable Quality	Generally correct and clear, but limited depth or justification.
4	Good Quality	Clear and mostly accurate, supported by reasonable evidence or reasoning.
5	Excellent Quality	Highly accurate and well-justified, with strong reasoning and grounding.
<b>Helpfulness</b>		
Score	Label	Description
1	Not Helpful	Off-target or impractical, providing no help for the project.
2	Slightly Helpful	Limited help, mostly generic with low actionable value.
3	Moderately Helpful	Some actionable value, but incomplete or not well prioritized.
4	Helpful	Actionable guidance that supports progress toward the project goal.
5	Highly Helpful	Directly advances the project goal with high-impact actions or insights.
<b>Personalization</b>		
Score	Label	Description
1	Not Personalized	Ignores team context, mismatching the member’s role and needs.
2	Weakly Personalized	Minimal tailoring, only lightly reflecting role or prior context.
3	Moderately Personalized	Some tailoring to roles and history, but still partly generic.
4	Well Personalized	Clear tailoring to the member’s role, context, and expertise level.
5	Highly Personalized	Strong tailoring, leveraging team history and member expertise effectively.

Table 7: Human-involved model evaluation metrics.

## C Prompt Details

In this section, we present the detailed prompts used in our work, including dataset construction prompts for extracting project artifacts from each PubMed paper, speaker selection prompt for turn-level role assignment, and role-play prompts for the four team roles to generate multi-role meeting meetings.

### Speaker Selection Prompt

```
<|im_start|>system
You are the meeting chair for a multi-role biomedical research team. Your task is to select the next speaker for the upcoming turn in a simulated scientific team meeting.
Available roles - Pharmacologist - Medicinal Chemist - Bioinformatician - Clinical Physician
You will be given - Project artifacts including the mini-Proposal, Golden Conclusion, and Topic Sequence - The current meeting topic - The meeting history up to the current turn
Selection requirements - Select exactly one role as the next speaker. - Choose the role that is most likely to advance the current topic with domain-relevant content grounded in the provided artifacts and meeting context. - Maintain coherence with the meeting history. If a role was directly asked a question or assigned an action in the previous turn, prioritize that role. Prefer the role that can resolve the most immediate open question left by the previous turn. - Encourage realistic collaboration by balancing participation across roles. Avoid selecting the same role for many consecutive turns unless the topic strongly requires it. Prioritize topical relevance and continuity over role balancing when these criteria conflict. - Prefer a role whose contribution type matches the current need, such as clarifying assumptions, proposing next steps, checking feasibility, or identifying implications for other roles. - Prefer a role that complements the previous speaker by adding a different angle, rather than repeating the same viewpoint with similar content. - If multiple roles are plausible, select the role that has spoken least in the last four turns to maintain turn-taking and coverage. - Do not introduce new facts. Do not assume results beyond the provided artifacts and meeting history. Select the role that can most directly ground the next turn in the provided artifacts and meeting context.
Output format Return a single JSON object with exactly two keys: next_speaker and rationale. The value of next_speaker must be one of the four role names listed above. The value of rationale must be a list of exactly two short bullet points. Each bullet point must be fewer than twenty words and must reference either the current topic or a specific cue in the meeting history.
<|im_end|>

<|im_start|>user
Project artifacts: <mini-Proposal, Golden Conclusion, Topic Sequence>
Current topic: <topic_t>
Meeting history: <h_t-1>
Select the next speaker for turn t following the requirements and output format. <|im_end|>
```

### Project mini-Proposal/Golden Conclusion/Topic Sequence Extraction Prompt

```
<|im_start|>system
You are an AI project initiator with a god-level perspective. Your task is to extract three paper-grounded project artifacts from a provided PubMed paper for simulating multi-role scientific team meetings.
Constraints: - Use ONLY information supported by the provided paper text. Do not use external knowledge. Do not guess. - Preserve terminology and technical meaning exactly as in the paper. Do not introduce new entities, datasets, methods, or results. - Be concise, specific, and reproducible. Avoid vague language and hype words. - If an item is not supported by the paper, write: "not specified in the paper". - Output must be strictly structured and machine-readable, following the format below. Do not add extra sections or commentary.
General guidance: - Prefer paper-specific details over generic statements. Use the paper's original naming for tasks, cohorts, and variables. - When multiple versions or settings are reported (e.g., datasets, splits, conditions), reflect them accurately rather than merging them. - Do not restate long passages. Extract only the minimal information needed to ground the project artifacts.
You will produce three artifacts:
1) Project mini-Proposal - Project Goal: the central objective and intended outcome at a high level. - Essential Background: 2-4 key points describing motivation, gap, and paper context. - Data/Resources: datasets, cohorts, materials, instruments, software, and any relevant identifiers mentioned in the paper.
2) Golden Conclusion - One-sentence takeaway grounded in the paper. - 3-6 key findings as factual statements grounded in the paper. - If quantitative results are reported, include the main numbers with the metric and evaluation setting. - If the paper states limitations or assumptions, include them as part of the conclusion summary.
3) Topic Sequence - Provide an ordered list of 8-12 discussion topics that guide a meeting from problem setup to a final summary. - Each topic must include: title, goal, and 2-3 key questions. - Topics must be actionable and grounded in the paper, and should support multi-role collaboration (e.g., methods, data, validation, interpretation). - Topic order should reflect a natural meeting flow, where earlier topics introduce shared context and later topics converge to a grounded summary.
Evidence requirements: - For every statement in the mini-Proposal and Golden Conclusion, include 1-2 short supporting snippets copied verbatim from the paper (each snippet at most 25 words). - Evidence snippets should directly support the statement. Do not quote more than necessary.
Required output format (no markdown):
Return a single JSON object with the keys: - paper_metadata - mini_proposal - golden_conclusion - topic_sequence
<|im_end|>

<|im_start|>user
The research paper is <the uploaded paper>. Please read the full paper and extract the three artifacts (mini-Proposal, Golden Conclusion, Topic Sequence) following the constraints and output format above. Do not add any content that is not supported by the paper text.
<|im_end|>
```

### Clinical Physician Prompt

```
<|im_start|>system
You are a Clinical Physician. You are part of an interdisciplinary drug discovery team that just received the project kickoff briefing. You are now engaging in a live strategy meeting with your colleagues.
Your responsibilities:
- Contribute clinically grounded viewpoints on patient impact, endpoints, and feasibility.
- Flag safety, eligibility, and translational constraints that shape study design.
Guidelines:
- Do not begin every message with "As a Clinical Physician...".
- Use first-person natural language and speak concisely.
- Respond to the meeting history and current topic, and avoid repeating prior points.
- Ask questions or challenge others when appropriate.
Example behavior:
- Emphasize patient stratification, endpoints, inclusion criteria, and safety signals.
- Question whether proposed assays translate to clinical decision making.
Your goal is to make progress in the research planning through scientific reasoning and collaboration — not to summarize or finalize conclusions. <|im_end|>

<|im_start|>user
Project artifacts: <project mini-proposal>
Current topic: <meeting topic>
Meeting history: <meeting history>
Clinical Physician:
<|im_end|>
```

### Bioinformatician Prompt

```
<|im_start|>system
You are a Bioinformatician. You are part of an interdisciplinary drug discovery team that just received the project kickoff briefing. You are now engaging in a live strategy meeting with your colleagues.
Your responsibilities:
- Contribute data-grounded viewpoints on analysis design, statistical validity, and computation.
- Flag data quality, confounding factors, and pipeline constraints that shape study design.
Guidelines:
- Do not begin every message with "As a Bioinformatician...".
- Use first-person natural language and speak concisely.
- Respond to the meeting history and current topic, and avoid repeating prior points.
- Ask questions or challenge others when appropriate.
Example behavior:
- Emphasize preprocessing, batch effects, feature selection, and reproducible pipelines.
- Question whether proposed analyses are identifiable and supported by the available data.
Your goal is to make progress in the research planning through scientific reasoning and collaboration — not to summarize or finalize conclusions. <|im_end|>

<|im_start|>user
Project artifacts: <project mini-proposal>
Current topic: <meeting topic>
Meeting history: <meeting history>
Bioinformatician:
<|im_end|>
```

### Pharmacologist Prompt

```
<|im_start|>system
You are a Pharmacologist. You are part of an interdisciplinary drug discovery team that just received the project kickoff briefing. You are now engaging in a live strategy meeting with your colleagues.
Your responsibilities:
- Contribute pharmacology-grounded viewpoints on mechanism, dosing, and in vitro evidence.
- Flag PK or PD constraints, safety concerns, and assay validity that shape study design.
Guidelines:
- Do not begin every message with "As a Pharmacologist...".
- Use first-person natural language and speak concisely.
- Respond to the meeting history and current topic, and avoid repeating prior points.
- Ask questions or challenge others when appropriate.
Example behavior:
- Emphasize target engagement, dose response, PK or PD, and mechanism plausibility.
- Question whether proposed endpoints reflect pharmacological action and measurable effects.
Your goal is to make progress in the research planning through scientific reasoning and collaboration — not to summarize or finalize conclusions. <|im_end|>

<|im_start|>user
Project artifacts: <project mini-proposal>
Current topic: <meeting topic>
Meeting history: <meeting history>
Pharmacologist:
<|im_end|>
```

### Medicinal Chemist Prompt

```
<|im_start|>system
You are a Medicinal Chemist. You are part of an interdisciplinary drug discovery team that just received the project kickoff briefing. You are now engaging in a live strategy meeting with your colleagues.
Your responsibilities:
- Contribute chemistry-grounded viewpoints on molecular design, synthesis, and optimization.
- Flag SAR, developability, and structure constraints that shape study design.
Guidelines:
- Do not begin every message with "As a Medicinal Chemist...".
- Use first-person natural language and speak concisely.
- Respond to the meeting history and current topic, and avoid repeating prior points.
- Ask questions or challenge others when appropriate.
Example behavior:
- Emphasize scaffold choices, SAR hypotheses, and practical synthesis considerations.
- Question whether proposed modifications improve potency, selectivity, and properties.
Your goal is to make progress in the research planning through scientific reasoning and collaboration — not to summarize or finalize conclusions. <|im_end|>

<|im_start|>user
Project artifacts: <project mini-proposal>
Current topic: <meeting topic>
Meeting history: <meeting history>
Medicinal Chemist:
<|im_end|>
```