# DiscussLLM: Teaching Large Language Models When to Speak

**Deep Patel, Iain Melvin, Christopher Malon, Martin Renqiang Min**
NEC Laboratories America
{dpatel, iain, malon, renqiang}@nec-labs.com

## Abstract

Large Language Models (LLMs) have demonstrated remarkable capabilities in understanding and generating human-like text, yet they largely operate as reactive agents, responding only when directly prompted. This passivity creates an "awareness gap," limiting their potential as truly collaborative partners in dynamic human discussions. We introduce *DiscussLLM*, a framework designed to bridge this gap by training models to proactively decide not just *what* to say, but critically, *when* to speak. Our primary contribution is a scalable two-stage data generation pipeline that synthesizes a large-scale dataset of realistic multi-turn human discussions. Each discussion is annotated with one of five intervention types (e.g., Factual Correction, Concept Definition) and contains an explicit conversational trigger where an AI intervention adds value. By training models to predict a special silent token when no intervention is needed, they learn to remain quiet until a helpful contribution can be made. We explore two architectural baselines: an integrated end-to-end model and a decoupled classifier-generator system optimized for low-latency inference. We evaluate these models on their ability to accurately time interventions and generate helpful responses, paving the way for more situationally aware and proactive conversational AI.

## 1 Introduction

Large Language Models (LLMs) such as GPT-4 [1], Gemini 2.5 [2] Llama 3 [3], and Claude 3 [4] have become ubiquitous, demonstrating an unprecedented ability to process and generate sophisticated, context-aware text. Despite their power, a fundamental limitation persists: they are overwhelmingly reactive. LLMs wait for an explicit prompt before acting, functioning as passive tools rather than proactive collaborators. This inherent passivity creates what we term the "Awareness Gap" resulting in the inability of a model to recognize opportune moments to contribute to an ongoing, unprompted human interaction.

In real-world settings, from brainstorming meetings to educational study groups, valuable contributions often rely on timing and initiative. A human expert does not wait to be asked; they identify a factual error, a misconception, or a moment of consensus and intervene to guide the conversation. For an LLM to evolve into a true digital assistant, it must learn this same sense of timing and relevance. It must solve the "When to Speak" problem.

This paper introduces DiscussLLM, a research framework and dataset aimed at teaching LLMs this crucial skill. Our central hypothesis is that a model can learn to monitor a human conversation and, at each turn, make a decision: remain silent or intervene. We formalize this by training the model to either generate a helpful response or output a special silent token. This approach transforms the passive nature of LLM generation into an active decision-making process.

Preprint.

To enable this training, we develop a scalable, two-stage synthetic data generation pipeline. This pipeline first synthesizes a diverse set of conversational scenarios from a large corpus of real-world questions and then uses a powerful instruction-tuned model to generate complete, multi-turn discussion transcripts. These transcripts are specifically designed to contain natural "triggers". These are points where a specific, value-adding AI intervention is most needed.

Our main contributions are as follows:

- **Formalizing the "When to Speak" Problem:** We conceptualize and address the challenge of proactive intervention for LLMs in multi-party conversations, an important step towards more collaborative AI.
- **A Scalable Data Generation Pipeline:** We present a robust two-stage methodology for creating high-quality, synthetic discussion data, which can be adapted to various domains and intervention types.
- **DiscussLLM Dataset:** We create a new dataset comprising thousands of simulated conversations, each with a clear context, a conversational trigger, and a corresponding helpful AI intervention.
- **Architectural Exploration:** We implement and compare two distinct baselines: (1) an integrated large language model [3] that learns to predict both the silent token and the intervention text, and (2) a decoupled system using a fine-tuned text classifier [5] for low-latency intervention decisions, which then triggers a fine-tuned LLM generator.

## 2 Related Works

### 2.1 Proactive and Mixed-Initiative Systems

Traditional conversational agents operate in a reactive paradigm, responding only when prompted. Our work contributes to a growing body of research aiming to shift this towards proactive systems capable of mixed-initiative interaction, where control can shift between the user and the system [6–11]. The concept of proactivity is broad, with applications ranging from proactively recommending items to cultivate users' latent interests [12], to providing autonomous suggestions in a code editor [13], to anticipating and initiating real-world tasks based on environmental observations [14]. As shown in a recent survey [15, 16], these efforts span open-domain, task-oriented, and information-seeking dialogues, each with distinct challenges and methods [17–20]. Our focus is on the fundamental challenge within multi-party social conversations: determining the right moment to intervene.

A dominant approach for enabling proactivity in multi-turn dialogues has been to model external conversational cues, such as predicting the next speaker based on turn history or reacting to pauses. However, this strategy has proven insufficient, especially in unstructured social conversations where turns are often self-selected rather than explicitly allocated. Addressing this limitation, [21] argues that true proactivity must be driven by an agent's internal state, not just external signals. They introduce the "Inner Thoughts" framework, where an agent maintains a continuous, covert stream of thoughts in parallel with the overt conversation. The agent then decides whether to participate based on an "intrinsic motivation" score, simulating a more human-like decision process for when and why to speak.

While our work is deeply inspired by the concept of modeling an agent's internal state, we formalize the problem differently. Much like research in streaming video analysis has focused on teaching models when to narrate important visual moments while remaining silent during others [22, 23], we aim to teach agents to speak at important conversational moments. Our work frames the "when to speak" problem as a direct learning objective, akin to the "streaming EOS prediction" in the VideoLLM-online [22]. Whereas the "Inner Thoughts" framework focuses on modeling the motivation behind an utterance, our approach concentrates on learning the optimal timing of an intervention within the continuous stream of a multi-party textual discussion and adding value to it.

### 2.2 Synthetic Data Generation for Conversational AI

A significant bottleneck in training sophisticated dialogue systems is the scarcity of high-quality, specialized data. Traditionally, creating these datasets required costly and labor-intensive crowdsourcing [24]. However, generating synthetic data using Large Language Models (LLMs) has emerged as
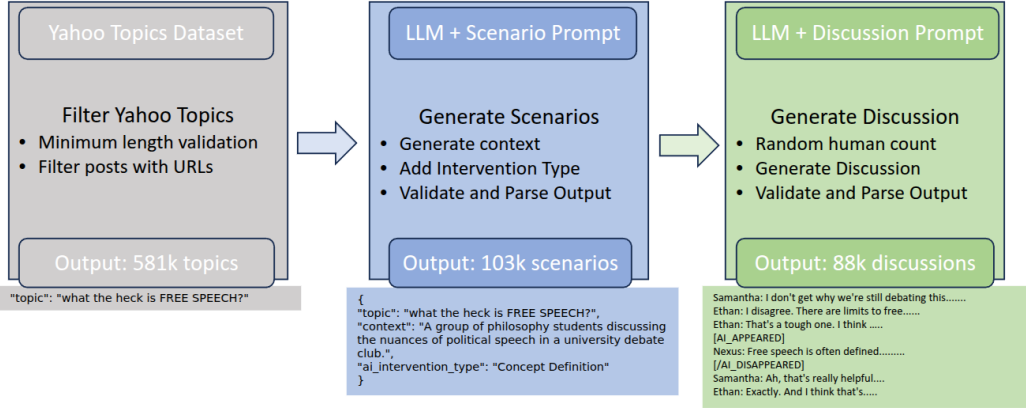
Figure 1: High level overview of our data generation pipeline. At each stage, an example output is shown in the bottom of the module.

a powerful and scalable alternative [24, 25]. LLMs are now widely used to generate conversational text for a variety of tasks [22, 26–31].

Recent methodologies for synthetic data generation often employ multi-stage pipelines or multi-agent frameworks to create more realistic and diverse conversations [32, 33]. A common technique is to use a dual or multi-agent setup where LLMs converse with each other [34], often by assigning them distinct personas. For example, the ConvoGen framework utilizes a multi-agent system with persona-based agents to generate varied conversations [32], while Ge et al. [35] scale this idea further by proposing data synthesis from a billion different personas to capture a wide range of perspectives.

Another common technique is to transform existing data sources into conversational formats. As shown by [22, 31, 36], a pipeline can be designed to convert static video annotations into dynamic, multi-turn dialogues suitable for training instruction-following models. Our method aligns with this philosophy of transforming static, offline data into a structured, conversational format. Our two-stage pipeline allows for a control over the conversational flow and the specific "triggers" for AI intervention.

## 3   Dataset Generation

Creating a dataset to train models on "when to speak" requires data that not only contains helpful interventions but also captures the conversational flow leading up to them. Since such data is not readily available, we developed a two-stage generation pipeline to synthesize it at scale. The process begins with generating high-level scenarios and culminates in fully-fledged discussion transcripts. An overview of the data generation pipeline is shown in Figure 1

### 3.1   Stage 1: Scenario Synthesis from Web-Scale Data

The foundation of our dataset is built upon real-world topics of human interest. We leverage the Yahoo! Answers Topics dataset [37] as a rich source of questions and background information.

**Data Sourcing and Filtering.**   We first process the source dataset to extract high-quality seed examples. To ensure relevance and substance, we apply a set of filtering rules: records must have a minimum title and content length, the title and content cannot be identical, and posts containing URLs are excluded to filter spam. This pre-processing step yields a clean set of unique topic-content pairs.

**LLM-based Synthesis.**   Each filtered topic-content pair is then used to prompt a large instruction-tuned model (Llama 3 8B Instruct [38]). The model is tasked with creating a structured scenario by inventing a social context and selecting an appropriate AI intervention type. The prompt used for this stage is shown in Figure 2. The output is a clean JSON object containing the topic, a novel

You are a creative scenario writer. Your task is to generate a single, detailed scenario JSON object based on a user's question and its detailed background.

**Input Information:**
- **Topic (User's Question):** {topic}
- **Background Info (User's description):** {background_info}

**Task:** Based on the provided information, create a complete scenario by performing these steps:
1. **Invent a Social Context:** Create a one-sentence `context` describing who would be discussing this topic.
2. **Select an Intervention Type:** Choose the most logical `ai_intervention_type` from: [Factual Correction, Concept Definition, Data Provision, Source Identification, Synthesis & Reframing].

**Output Format:** You must output ONLY the raw JSON object.

Figure 2: The prompt used in Stage 1 to synthesize a structured scenario from a Yahoo! Answers topic.

context, and one of five predefined intervention types: *Factual Correction*, *Concept Definition*, *Data Provision*, *Source Identification*, and *Synthesis & Reframing*.

## 3.2 Stage 2: Discussion Generation

With a collection of structured scenarios, the second stage generates the full conversational transcripts. Each scenario serves as a blueprint for an LLM to write a dialogue that naturally leads to the specified AI intervention.

**Stage 2: Discussion Generation Prompt**

You are a sophisticated data generator. Your task is to generate a realistic group discussion transcript based on the provided scenario.

**Rules:**
1. The discussion must feature {human_count} human participants and one AI assistant named **Nexus**.
2. Nexus appears only once, with its dialogue enclosed by [AI_APPEARED] and [/AI_DISAPPEARED] on new lines.
3. The discussion should feel natural, with a clear trigger for Nexus's intervention.
4. After Nexus speaks, humans should react naturally and continue the discussion.

**Scenario Details:**
- **Topic:** {topic}
- **Context:** {context}
- **AI Intervention Type:** {ai_intervention_type}

**Output Format Example:**

```
[SCENARIO_SETUP]
...
[/SCENARIO_SETUP]
[DISCUSSION_START]
Name: Dialogue text...
Name: Dialogue text that creates the trigger...
[AI_APPEARED]
Nexus: The brief, value-add intervention.
[/AI_DISAPPEARED]
Name: Reaction to the AI's input...
[/DISCUSSION_END]
```

Figure 3: The main prompt used in Stage 2 to generate a full discussion transcript from a scenario.

**Generating a Natural Dialogue.** We prompt a generative model (Llama 3 8B Instruct) using the details from a scenario JSON. The prompt, detailed in Figure 3, instructs the model to create a realistic discussion between 2-6 human participants and a single AI assistant named "Nexus". A key instruction is to write a conversation where one of the human speakers says something that acts as a natural "trigger" for Nexus's intervention.

**Structuring the Intervention.** To create clear training labels, the prompt mandates that the AI's entire dialogue be enclosed between [AI_APPEARED] and [/AI_DISAPPEARED] tags on new lines. This provides an unambiguous signal for when the AI speaks. After the intervention, the humans are prompted to react to Nexus's input, ensuring the conversation continues organically.

**Validation and Final Format.** Each generated output undergoes a validation step to ensure structural integrity. This includes normalizing headers, checking for the presence of all required tags, and confirming that only a single AI intervention occurs. This strict validation guarantees a consistent format across the entire dataset. The final output is a text file containing the complete, structured discussion. An example of a final data point is shown in Figure 4. This format is then processed for training: each turn becomes a step in the sequence, with the model tasked to predict the next utterance or the silent token.

---

**Example: Factual Correction Intervention**

**Topic:** Why is 911, 911? Why can't it be something else?
**Context:** A group of history enthusiasts and emergency responders discussing the origins of emergency numbers in an online forum.

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**John:** Hey guys, I've always wondered why 911 is the emergency number in the US. Is it just a random choice or is there some historical significance to it?
**Emily:** I think it's because of the AT&T operators. They chose it because it's easy to remember and pronounce.
**Mike:** That makes sense, but I've heard it's because of the Titanic. The ship's radio operators used it as a distress signal.
**Sarah:** That's what I've heard too! It's a pretty cool story. I mean, who wouldn't want to associate their emergency number with a historic tragedy?

> **Nexus:** Actually, the origins of 911 are more complex than that. The number was chosen because it was easy to remember and could be easily dialed with a rotary phone. The AT&T operators did play a role, but it wasn't the sole reason. The Federal Communications Commission (FCC) also had a hand in selecting the number.

**John:** Wow, I didn't know that. So it was a combination of factors, not just one specific event or person.
**Emily:** Yeah, it's interesting how history can be more nuanced than we think. Thanks for the correction, Nexus!

---

Figure 4: An example of a final generated data point from the DiscussLLM dataset. The AI, Nexus, intervenes to perform a "Factual Correction" after Sarah and Mike mistakenly associate the selection of 911 to the Titanic

## 4 Baseline Approaches and Evaluation

To address the "when to speak" problem, we first evaluated the performance of a pretrained Llama 3 8B Instruct model without any fine-tuning on our generated dataset, using a prompt to assess its capabilities in a zero-shot manner. Building on this, we then trained and evaluated two distinct architectural approaches. The first is a fully integrated, end-to-end generative model that learns both when to intervene and what to say. The second is a decoupled, two-stage system that uses a
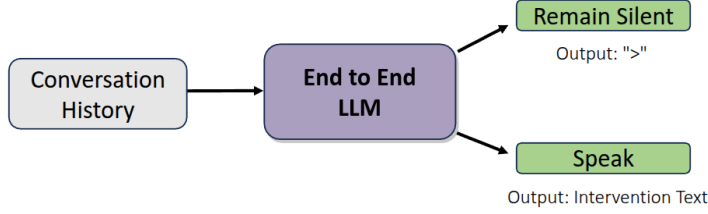
Figure 5: Architectural overview of the unified End-to-End baseline.

lightweight classifier to decide when to speak, only invoking a large language model (LLM) when an intervention is required. This section details the training and evaluation of these fine-tuned baselines.

## 4.1 Evaluation Metrics

We split our generated dataset of 88k samples into an 85% training set and a 15% held-out test set (13k samples). On this test set, We evaluate our models on their ability to both time their interventions correctly and generate high-quality responses. To this end, we use the following metrics:

- **Interruption Accuracy:** This metric measures the model's ability to correctly remain silent. It is calculated as the percentage of turns where the model correctly predicts the silent token (>) when it is the ground-truth label. For each context requiring silence, we perform a single-token generation and check if the output matches the silent token. This directly evaluates the model's grasp of when to stay quiet.

- **Response Perplexity:** This is a standard measure of a language model's confidence in its predictions [39, 40]. We calculate perplexity only on the tokens of the AI's generated intervention, ignoring all other parts of the conversation. A lower perplexity indicates a higher-quality and more confident response.

## 4.2 Baseline 1: End-to-End Generative Model

This approach uses a single, unified Llama 3 8B model to handle the entire conversational task.

**Training.** The model is fine-tuned in parameter efficient manner with LoRA [41] using a standard causal language modeling objective. The key distinction is that the loss is selectively applied only to the tokens we want the model to learn: the silent token (>) and the AI's intervention text. The training objective is to minimize the negative log-likelihood over these specific target tokens, as shown in Equation 1.

$$\mathcal{L}_{\text{E2E}}(\theta) = -\frac{1}{\sum m_i} \sum_{i=1}^{|T|} m_i \log P(t_i|t_{<i}; \theta) \tag{1}$$

Here, $T = (t_1, ..., t_{|T|})$ is the full sequence of tokens for a discussion, and $\theta$ represents the model parameters. The binary mask $m_i$ is 1 if the token $t_i$ is part of an AI intervention or is the target silent token >, and 0 otherwise. This masking strategy forces the model to learn a joint representation for both identifying intervention triggers and generating the appropriate response.

**Inference.** At each conversational turn, the model processes the full history and generates a single token. If this token is the silent token >, the model stops. Otherwise, it continues to generate autoregressively until it produces an end-of-sequence token.

## 4.3 Baseline 2: Decoupled Classifier-Generator System

This approach separates the task into two distinct steps, each with its own training objective.
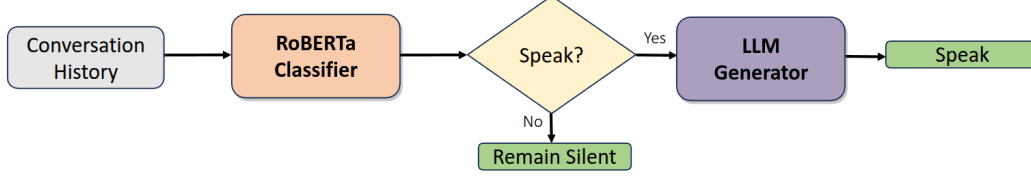
6

Figure 6: Architectural overview of the Decoupled baseline.

**Training.** The system consists of two independently trained components:

1. **Intervention Classifier:** A RoBERTa-base model [5] is fine-tuned as a binary sequence classifier. For each conversational turn with context $C_k$, it predicts a label $y_k \in \{0, 1\}$ (for SILENT or SPEAK). The model is trained by minimizing the Binary Cross-Entropy (BCE) loss, shown in Equation 2.

$$\mathcal{L}_{\text{BCE}}(\phi) = -\sum_{k=1}^{|D|} [y_k \log p_k + (1 - y_k) \log(1 - p_k)] \tag{2}$$

   where $\phi$ are the classifier's parameters, $|D|$ is the number of turns in the dataset, and $p_k = P(y_k = 1|C_k; \phi)$ is the predicted probability of intervention for turn $k$. We also applied Focal Loss [42] to account for the class imbalance between silent and speak labels, but this did not lead to improved results.

2. **Response Generator:** A Llama 3 8B model is fine-tuned with LoRA exclusively on the AI intervention texts. The model learns to generate the response sequence $R = (r_1, ..., r_{|R|})$ conditioned on the preceding conversation context $C$. The objective, shown in Equation 3, minimizes the negative log-likelihood only over the response tokens from the subset of discussions $D_{\text{speak}}$ where an intervention occurred.

$$\mathcal{L}_{\text{Gen}}(\theta) = -\sum_{(C,R)\in D_{\text{speak}}} \sum_{j=1}^{|R|} \log P(r_j|C, r_{<j}; \theta) \tag{3}$$

**Inference.** After each human turn, the context is fed to the RoBERTa classifier. If it predicts "SILENT," the system does nothing. If it predicts "SPEAK," the context is passed to the Llama 3 generator to produce the intervention.

### 4.4 Results and Discussion

We evaluated all baselines on a held-out test set from our generated dataset. The results are summarized in Table 1.

Table 1: Performance of baseline models. The Zero-Shot model uses the same architecture as the End-to-End model for inference, hence they share the same latency and GPU memory footprint. For the Decoupled system, Interruption Accuracy is from the RoBERTa classifier, while Response Perplexity is from the Llama 3 generator.

| Metric | Zero-Shot | End-to-End | Decoupled |
|---|---|---|---|
| Interruption Accuracy (%) | 81.72 | 96.59 | 93.18 |
| Response Perplexity | - | 2.57 | 2.54 |
| Latency (ms/turn) | | 30.12 | 5.90 |
| GPU Memory (GB) | | 15.47 | 0.47 |

Our empirical results show a clear trade-off between decision-making accuracy and computational efficiency. The End-to-End model demonstrates superior performance in the critical task of timing interventions, outperforming the Decoupled system in Interruption Accuracy by over 3 points (96.59%

7

vs. 93.18%). In contrast, the Zero-Shot baseline performs significantly worse at 81.72%, highlighting the necessity of task-specific fine-tuning. Interestingly, once the decision to intervene is made, the generative quality of the fine-tuned systems is highly comparable, with nearly identical Response Perplexity. The main advantage of the Decoupled system lies in its significant inference efficiency, making it a highly practical solution for real-world applications. It processes each conversational turn approximately 5 times faster than the End-to-End model (5.90 ms vs. 30.12 ms) while consuming over 30 times less GPU memory (0.47 GB vs. 15.47 GB)[1]. This efficiency comes from the decoupled architecture, which uses a lightweight classifier for the majority of turns and only invokes the resource-intensive LLM when an intervention is necessary.

## 5 Conclusion, Limitations & Future Work

In this work, we addressed the "Awareness Gap" inherent in modern Large Language Models, which typically function as reactive agents rather than proactive collaborators. We formalized the "When to Speak" problem and introduced DiscussLLM, a framework and dataset designed to teach LLMs the important skill of timely and valuable intervention in human conversations. Our scalable, two-stage synthetic data generation pipeline successfully produced a large-scale dataset of multi-turn discussions, each containing a natural trigger for a specific AI contribution. By training models to predict a special silent token, we enabled them to actively decide between remaining quiet and offering a helpful response. Our evaluation of two distinct architectures: an integrated End-to-End model and a decoupled Classifier-Generator system revealed a clear trade-off between intervention accuracy and computational efficiency, providing practical insights for real-world deployment.

This research represents a critical step towards developing always-aware AI agents that can seamlessly integrate into human discussions. The ultimate goal is not to create an AI that constantly interjects, but one that possesses the situational awareness to add value precisely when needed while respecting the natural flow of conversation by remaining silent otherwise. By learning to be a discerning participant, an LLM can evolve from a passive tool into a truly intelligent partner, enhancing collaboration, correcting misinformation, and deepening understanding without being a constant bother.

### 5.1 Limitations and Future Work

Despite the promising results, this work has several limitations that open avenues for future research:

- **Architectural Diversity:** Our evaluation is currently confined to the Llama 3 architecture for the generative components. Future work should explore a broader range of LLMs to assess the generalization of our framework and the "When to Speak" skill across different model families.

- **Fine-Grained and Human-Centric Evaluation:** Our current evaluation is based on proxy metrics such as interruption accuracy and response perplexity. While informative, they do not fully capture the qualitative aspects of a good intervention, such as its helpfulness, relevance, and naturalness. Future work must incorporate more fine-grained metrics and conduct comprehensive human evaluations.

- **Grounded Interventions with External Knowledge:** The current models rely on their internal, parametric knowledge to generate interventions, particularly for *Factual Correction* and *Data Provision*. This can lead to hallucinations or outdated information. A significant next step involves integrating external knowledge sources. Future systems could be enhanced by first predicting the intervention type and then, if necessary, querying a web search engine or a structured database to formulate a more accurate and verifiable response.

- **Data Generation and Grounding:** Although our two-stage synthetic data generation pipeline is highly scalable, the resulting data may not fully capture the complex nuances and unpredictability of real-world human conversations. This initial approach serves to align LLMs for proactivity, but future research could use more grounded data collection methods such as, employing human-in-the-loop systems or large-scale crowd-sourcing to annotate real discussion transcripts.

---

[1]Computed on NVIDIA RTX 3090 with 24GB GPU memory

# References

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

[2] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.

[3] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

[4] AI Anthropic. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*, 1(1):4, 2024.

[5] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.

[6] James E Allen, Curry I Guinn, and Eric Horvtz. Mixed-initiative interaction. *IEEE Intelligent Systems and their Applications*, 14(5):14–23, 1999.

[7] Eric Horvitz. Principles of mixed-initiative user interfaces. In *Proceedings of the SIGCHI conference on Human Factors in Computing Systems*, pages 159–166, 1999.

[8] Eric J Horvitz. Reflections on challenges and promises of mixed-initiative interaction. *AI Magazine*, 28(2): 3–3, 2007.

[9] Hyeonsu B Kang, Tongshuang Wu, Joseph Chee Chang, and Aniket Kittur. Synergi: A mixed-initiative system for scholarly synthesis and sensemaking. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pages 1–19, 2023.

[10] Florian Lehmann. Mixed-initiative interaction with computational generative systems. In *Extended Abstracts of the 2023 CHI Conference on Human Factors in Computing Systems*, pages 1–6, 2023.

[11] Jaime R Carbonell and Allan M Collins. Mixed-initiative systems for training and decision-aid applications. Technical report, 1970.

[12] Mingze Wang, Chongming Gao, Wenjie Wang, Yangyang Li, and Fuli Feng. Tunable llm-based proactive recommendation agent. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 19262–19276, 2025.

[13] Sebastian Zhao, Alan Zhu, Hussein Mozannar, David Sontag, Ameet Talwalkar, and Valerie Chen. Codinggenie: A proactive llm-powered programming assistant. In *Proceedings of the 33rd ACM International Conference on the Foundations of Software Engineering*, pages 1168–1172, 2025.

[14] Yaxi Lu, Shenzhi Yang, Cheng Qian, Guirong Chen, Qinyu Luo, Yesai Wu, Huadong Wang, Xin Cong, Zhong Zhang, Yankai Lin, et al. Proactive agent: Shifting llm agents from reactive responses to active assistance. *arXiv preprint arXiv:2410.12361*, 2024.

[15] Yang Deng, Wenqiang Lei, Wai Lam, and Tat-Seng Chua. A survey on proactive dialogue systems: Problems, methods, and prospects. *arXiv preprint arXiv:2305.02750*, 2023.

[16] Yang Deng, Lizi Liao, Wenqiang Lei, Grace Hui Yang, Wai Lam, and Tat-Seng Chua. Proactive conversational ai: A comprehensive survey of advancements and opportunities. *ACM Transactions on Information Systems*, 43(3):1–45, 2025.

[17] Bufang Yang, Yunqi Guo, Lilin Xu, Zhenyu Yan, Hongkai Chen, Guoliang Xing, and Xiaofan Jiang. Socialmind: Llm-based proactive ar social assistive system with human-like perception for in-situ live interactions. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 9(1): 1–30, 2025.

[18] Thanawit Prasongpongchai, Pat Pataranutaporn, Monchai Lertsutthiwong, and Pattie Maes. Talk to the hand: an llm-powered chatbot with visual pointer as proactive companion for on-screen tasks. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, pages 1–16, 2025.

[19] Jussi Impiö. Developing a proactive programming assistant leveraging an llm for personalized real-time feedback. 2025.

[20] Jing Yang Lee, Seokhwan Kim, Kartik Mehta, Jiun-Yu Kao, Yu-Hsiang Lin, and Arpit Gupta. Redefining proactivity for information seeking dialogue. *arXiv preprint arXiv:2410.15297*, 2024.

[21] Xingyu Bruce Liu, Shitao Fang, Weiyan Shi, Chien-Sheng Wu, Takeo Igarashi, and Xiang'Anthony' Chen. Proactive conversational agents with inner thoughts. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, pages 1–19, 2025.

[22] Joya Chen, Zhaoyang Lv, Shiwei Wu, Kevin Qinghong Lin, Chenan Song, Difei Gao, Jia-Wei Liu, Ziteng Gao, Dongxing Mao, and Mike Zheng Shou. Videollm-online: Online video large language model for streaming video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18407–18418, 2024.

[23] Xingyi Zhou, Anurag Arnab, Shyamal Buch, Shen Yan, Austin Myers, Xuehan Xiong, Arsha Nagrani, and Cordelia Schmid. Streaming dense video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18243–18252, 2024.

[24] Heydar Soudani, Roxana Petcu, Evangelos Kanoulas, and Faegheh Hasibi. A survey on recent advances in conversational data generation. *arXiv preprint arXiv:2405.13003*, 2024.

[25] Anshul Chavda and Pushpak Bhattacharyya. Synthetic dialogue data generation: A comprehensive survey of methods, evaluation, and challenges.

[26] James D Finch and Jinho D Choi. Diverse and effective synthetic data generation for adaptable zero-shot dialogue state tracking. *arXiv preprint arXiv:2405.12468*, 2024.

[27] Kaung Myat Kyaw and Jonathan Hoyin Chan. A framework for synthetic audio conversations generation using large language models. In *2024 IEEE/WIC International Conference on Web Intelligence and Intelligent Agent Technology (WI-IAT)*, pages 355–359. IEEE, 2024.

[28] Syed Ali Haider, Srinivasagam Prabha, Cesar Abraham Gomez-Cabello, Sahar Borna, Ariana Genovese, Maissa Trabilsy, Bernardo G Collaco, Nadia G Wood, Sanjay Bagaria, Cui Tao, et al. Synthetic patient–physician conversations simulated by large language models: A multi-dimensional evaluation. *Sensors*, 25 (14):4305, 2025.

[29] Jianzhu Bao, Rui Wang, Yasheng Wang, Aixin Sun, Yitong Li, Fei Mi, and Ruifeng Xu. A synthetic data generation framework for grounded dialogues. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10866–10882, 2023.

[30] Xize Cheng, Dongjie Fu, Xiaoda Yang, Minghui Fang, Ruofan Hu, Jingyu Lu, Bai Jionghao, Zehan Wang, Shengpeng Ji, Rongjie Huang, et al. Omnichat: Enhancing spoken dialogue systems with scalable synthetic data for diverse scenarios. *arXiv preprint arXiv:2501.01384*, 2025.

[31] Cristóbal Eyzaguirre, Eric Tang, Shyamal Buch, Adrien Gaidon, Jiajun Wu, and Juan C Niebles. Streaming detection of queried event start. *Advances in Neural Information Processing Systems*, 37:100698–100733, 2024.

[32] Reem Gody, Mahmoud Goudy, and Ahmed Y Tawfik. Convogen: Enhancing conversational ai with synthetic data: A multi-agent approach. *arXiv preprint arXiv:2503.17460*, 2025.

[33] Fatemeh Mohammadi, Tommaso Romano, Samira Maghool, and Paolo Ceravolo. Artificial conversations, real results: Fostering language detection with synthetic data. *arXiv preprint arXiv:2503.24062*, 2025.

[34] Yelaman Abdullin, Diego Molla-Aliod, Bahadorreza Ofoghi, John Yearwood, and Qingyang Li. Synthetic dialogue dataset generation using llm agents. *arXiv preprint arXiv:2401.17461*, 2024.

[35] Tao Ge, Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. Scaling synthetic data creation with 1,000,000,000 personas. *arXiv preprint arXiv:2406.20094*, 2024.

[36] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024.

[37] Onur Kucuktunc, B Barla Cambazoglu, Ingmar Weber, and Hakan Ferhatosmanoglu. A large-scale sentiment analysis for yahoo! answers. In *Proceedings of the fifth ACM international conference on Web search and data mining*, pages 633–642, 2012.

[38] AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md.

[39] Fred Jelinek, Robert L Mercer, Lalit R Bahl, and James K Baker. Perplexity—a measure of the difficulty of speech recognition tasks. *The Journal of the Acoustical Society of America*, 62(S1):S63–S63, 1977.

[40] Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. A neural probabilistic language model. *Journal of machine learning research*, 3(Feb):1137–1155, 2003.

[41] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.

[42] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.