# LLM-Enhanced Dialogue Management for Full-Duplex Spoken Dialogue Systems

Hao Zhang\*1, Weiwei Li<sup>2</sup>, Rilin Chen<sup>3</sup>, Vinay Kothapally<sup>1</sup>, Meng Yu<sup>1</sup>, Dong Yu<sup>1</sup>

<sup>1</sup>Tencent AI Lab, Bellevue, USA <sup>2</sup>Tencent AI Lab, Shenzhen, China <sup>3</sup>Tencent AI Lab, Beijing, China

aaronhzhang@global.tentcent.com

## **Abstract**

Achieving full-duplex communication in spoken dialogue systems (SDS) requires real-time coordination between listening, speaking, and thinking. This paper proposes a semantic voice activity detection (VAD) module as a dialogue manager (DM) to efficiently manage turn-taking in full-duplex SDS. Implemented as a lightweight (0.5B) LLM fine-tuned on fullduplex conversation data, the semantic VAD predicts four control tokens to regulate turn-switching and turn-keeping, distinguishing between intentional and unintentional barge-ins while detecting query completion for handling user pauses and hesitations. By processing input speech in short intervals, the semantic VAD enables real-time decision-making, while the core dialogue engine (CDE) is only activated for response generation, reducing computational overhead. This design allows independent DM optimization without retraining the CDE, balancing interaction accuracy and inference efficiency for scalable, nextgeneration full-duplex SDS.

**Index Terms**: full-duplex, spoken dialogue system, voice activity detection, dialog management, large language model

## 1. Introduction

Spoken dialogue systems (SDS) [1, 2, 3, 4] have advanced significantly with large language models (LLMs), enabling more natural and context-aware interactions [5, 6]. However, achieving full-duplex communication, where SDS can listen and speak simultaneously, remains challenging [7, 8, 9, 10]. Many SDSs still operate in a turn-based (half-duplex) or "pseudo" full-duplex manner [11, 12], leading to less fluid interactions due to lack of understanding of the user's state and intention. In contrast, human conversations involve seamless turn-taking [13]. Beyond turn-taking, full-duplex SDS must handle challenges like interfering speakers, user hesitations, and distinguishing between intentional and unintentional interruptions to improve naturalness and efficiency [14, 15, 16].

Achieving full-duplex communication in SDS has been widely studied. Early approaches relied on separate neural network modules for detection and classification tasks, limiting stability and robustness. Shin et al. [17] used speech event detection for real-time query updates but lacked broader interaction handling. Lin et al. [18] introduced Duplex Conversation, leveraging a multimodal model for user state detection and turn management. Recent methods embed interaction handling directly into LLMs, improving automation but increasing inference overhead. Wang et al. [19] proposed NeuralFSM, fine-tuning an LLM with a finite state machine for synchronized speaking and listening. Similarly, Zhang et al. [20] fine-

tuned an LLM with a time-division multiplexing strategy for duplex dialogue. VITA [21] incorporated state tokens and a duplex scheme for non-awakening and audio-interrupt interactions, while Moshi [22] integrated a base LLM with a smaller Transformer for real-time streaming predictions. Mai et al. [23] introduced RTTL-DG, a textless spoken dialogue model incorporating backchannels and laughter for natural turn-taking.

To balance conversational quality, stability, and inference efficiency, this paper proposes a semantic voice activity detection (VAD) as a dialogue manager (DM) for full-duplex SDS. While an acoustic VAD mitigates interference from background speakers using acoustic cues, the semantic VAD leverages LLM capabilities to detect user states and intentions based on semantic information. We implement the semantic VAD by fine-tuning a smaller (0.5B) LLM on carefully designed full-duplex text conversations. It predicts four control tokens-start-speaking, start-listening, continue-speaking, and continue-listening—to dynamically guide system interaction. Operating at short intervals, the DM enables real-time interaction management, distinguishing between intentional and unintentional barge-ins and detecting query completion to handle user pauses and hesitations. Our method improves upon prior work by balancing DM performance and computational efficiency: leveraging the semantic understanding of LLMs enable precise dialogue management, while a lightweight LLM cuts computational costs. Independent DM and CDE optimization further ensures scalability.

## 2. Full-duplex spoken dialogue system

## 2.1. Key challenges in full-duplex SDS

Achieving seamless full-duplex communication in SDS requires accurately detecting and differentiating user activities, states, and intentions. However, real-time turn-taking remains challenging due to issues such as interfering speakers, ambiguous pauses, and unintentional interruptions:

- Interfering Speakers: Background speech can cause ASR and dialogue management errors, leading to unintended activations or responses.
- User Pauses & Hesitations: Silence alone does not indicate query completion, often resulting in premature responses or unnecessary delays.
- Unintentional Interruptions: Acknowledgments, backchanneling, or speech directed at others may mistakenly halt SDS responses, disrupting conversational flow.

Addressing these challenges requires integrating both acoustic and semantic information for robust, context-aware full-duplex interaction.

<sup>\*</sup>Corresponding author

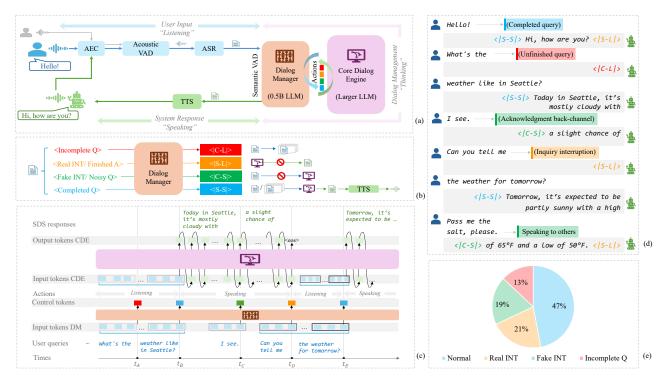


Figure 1: (a) The proposed system architecture; (b) Control tokens and their corresponding actions; (c) Interactions between the DM and CDE; (d) An example of a full-duplex conversation; and (e) Distribution of interaction scenarios in the generated full-duplex conversation dataset.

#### 2.2. Proposed full-duplex SDS design

The proposed SDS consists of six key modules: acoustic echo cancellation (AEC) [24, 25], acoustic VAD, automatic speech recognition (ASR), semantic VAD, an LLM serving as the core dialogue engine (CDE), and text-to-speech (TTS), as shown in Figure 1(a). User speech is processed sequentially, with the first four modules operating in short intervals to ensure real-time responsiveness. The CDE is activated only when needed, minimizing computational overhead.

We combine acoustic and semantic VAD to handle the previously mentioned challenges, enabling robust and seamless full-duplex interaction:

- Acoustic VAD: Our designed acoustic VAD is speaker- and distance-aware, leveraging distance information and speaker embeddings (if provided) to isolate target speaker's activity and mitigate interference from other speakers.
- Semantic VAD: Serving as the DM and implemented as a fine-tuned 0.5B LLM, it leverages the LLM's semantic understanding and contextual awareness to predict control tokens for dynamic turn-taking management. This enables effective handling of hesitations and accurate differentiation between real and fake interruptions.

This study focuses on introducing the semantic VAD. Rather than improving question-answering capabilities, our goal is to train the DM to follow structured rules and take appropriate actions in various full-duplex interaction scenarios through instruction tuning [26, 19]. This approach requires only a small amount of high-quality data and minimal fine-tuning, making it efficient and practical. Unlike NN-based DMs that rely on shallow features [17, 18], the fine-tuned LLM leverages semantic understanding for precise interaction manage-

ment. Compared to integrating DM directly into CDE (larger models) [19, 21, 20, 22], our method reduces inference costs by offloading frequent decisions to a lightweight LLM while ensuring scalability through independent DM and CDE optimization.

## 3. Semantic VAD for dialogue management

## 3.1. Interaction scenarios and action tokens

Most user-SDS interactions follow a standard questionanswering flow without interruptions. However, when interruptions occur, we prioritize user experience in our design by allowing only user-initiated interruptions.

The DM enables SDS to switch between speaking and listening modes by performing two key tasks, each associated with specific action tokens (shown in Figure 1(b)):

- User State Detection: Determines whether the user has finished speaking.
- Query Incomplete: The DM predicts < |Continue-Listening|> (< |C-L|>), allowing the system to keep listening and cache historical queries until a complete query is identified.
- Query Complete: The DM prompts the SDS to <|Start-Speaking|> (<|S-S|>). Once the response is finished, the SDS transitions to <|Start-Listening|> (<|S-L|>) again and awaits further input.
- User Intention Analysis: Determines whether a user bargein requires a response.
  - Intentional Interruption (Real INT): If the user intends to redirect the conversation, the SDS stops speaking and switches to <|Start-Listening|>(<|S-L|>). Ex-

#### Algorithm 1: Full-Duplex Data Prompts Generation

```
Input: Topic Pool T: Speaking Style Pool S:
            Hyperparameters: N_{conv} (number of
  conversations), P_{real} (probability of Real INT), P_{fake}
  (Fake INT), P_{incomplete} (incomplete Q)
  Output: Generated prompts for full-duplex
            conversation data
1 for i \leftarrow 1 to N_{conv} do
      Randomly select topic t and speaking style s;
2
3
       Set QA rounds n \sim \text{Uniform}(2, 12);
       Initialize conversation C with t and s;
4
       for j \leftarrow 1 to n do
5
           Randomly assign QA type with probability P_*:
            if QA round is Real INT (P_{real}) then
               Add real interruption scenario to C;
 7
           else if QA round is Fake INT(P_{fake}) then
 8
            Add fake interruption scenario to C;
10
11
              Add normal QA exchange to C;
          Randomly truncate user queries for incomplete
12
            queries (P_{incomplete});
       Generate a prompt by explaining control tokens
13
        placement and usage with examples;
       Append prompt to prompt set;
14
```

- amples include denial or discontent, further inquiries, and topic changes.
- Unintentional Interruption (Fake INT): If the interruption is non-disruptive, the SDS continues speaking with < | Continue-Speaking |> (< | C-S |>). Examples include affirmative acknowledgments, back-channeling, speech directed at someone else, and unrelated comments.

#### 3.2. Dialogue management

15 Return the set of generated prompts;

As shown in Figure 1(c), the DM processes two primary inputs: (1) user query tokens from ASR and (2) historical response tokens from the CDE. For clarity and due to space constraints, the historical response tokens are not explicitly depicted in the diagram as DM input. The DM regulates interaction flow by analyzing these inputs to determine the next action.

The DM serves as the primary decision-maker, continuously managing listening and speaking states, while the CDE focuses on generating high-quality responses when required. The CDE is activated only when the DM predicts < | S-S | >, prompting it to perform standard LLM inference using past user queries and system responses. By relying on the DM for frequent real-time decisions and selectively invoking the CDE, the system optimizes both computational efficiency and conversational quality, ensuring a scalable and effective full-duplex SDS.

## 3.3. Full-duplex conversation data

This section focuses on the core of the semantic VAD, detailing the data generation and the methodology for designing tailored prompts.

Creating high-quality training data for the dialogue manager (DM) requires annotated full-duplex conversations with the four control tokens, as shown in Figure 1(d). As no public

Table 1: Confusion matrix full-duplex SDS evaluation. "GT" denotes ground truth, and "Est" denotes estimation.

$GT \backslash Est$	<   C-L   >	< S-S >	< S-L >	< C-S >	Recall
<   C-L   >	926	74	0	0	0.926
< S-S >	11	989	0	0	0.989
< S-L >	1	0	999	0	0.999
<   C-S   >	0	0	0	1000	1.000
Precision	0.987	0.930	1.000	1.000	Accuracy:
F1 Score	0.956	0.959	0.999	1.000	0.9785

datasets meet these criteria, we used LLM API, Yuanbao [27], to generate diverse and natural full-duplex dialogues.

To generate prompts to create full duplex conversation data, we designed an algorithm as shown in Algorithm 1. The prompts were designed to instruct LLMs on how to generate full-duplex conversation data annotated with control tokens. Each prompt includes the following elements:

- **Objective**: Generate a conversation transcript where the user can interrupt the assistant.
- Assistant's Behavior: Specify the meanings of control tokens and corresponding actions.
- Output Format: Conversations are formatted as: {Round X (dialogue type); User: <user's query or barge-in>; Sys: <assistant's response with control tokens>}.
- Custom Instructions: Specify the number of QA rounds, topic, and speaking style. For example: "Please generate 5 rounds of conversation. The user initiates the conversation with a topic related to travel, speaking in a casual and conversational style. The assistant maintains a friendly tone and references the previous context."

The topic pool covers 200 common conversation themes (e.g., weather, travel, restaurants), while the speaking style pool includes 10 user personas. Topics and styles are randomly selected, with probabilities adjusted for realism.

## 3.4. Data preparation and model training

Data preparation consists of four stages: data generation, cleaning, post-processing, and augmentation. We generated full-duplex conversations using 20,000 prompts with Yuanbao [27], but found that not all generated data align fully with the intended structure due to limited command-following abilities. To address this, we applied strict data cleaning, filtering malformed outputs, and incorrect control tokens, resulting in a retention rate of 60%. Given this, we also adopted an alternative strategy: first generating standard QA dialogues, which LLMs handle more effectively, and then introduce various interaction patterns via controlled post-processing. The final training dataset combines conversations from both strategies.

We further augmented the data by modifying the punctuation at the end of user queries to better match the ASR output. The final dataset includes 11,990 conversations with 80,338 dialogue rounds covering diverse interaction scenarios (Figure 1(e)). Furthermore, we examine extreme cases with a high proportion (over 50%) or exclusively real/fake interruptions or incomplete queries, which led to overfitting and degraded performance, reinforcing the need for balanced training data. To mitigate potential degradation of the base LLM's capabilities,

Table 2: Comparison with related studies.

Tasks & Actions		Precision		Recall			F1				
DuplexConv	RTTL-DG	SemanticVAD	[18]	[23]	SemanticVAD	[18]	[23]	SemanticVAD	[18]	[23]	SemanticVAD
[18]	[23]	(proposed)	[10]	[20]	Semantic vi is	[10]	[20]		[10]	[20]	Delinantie (112
User state detection	Remain	Continue listening	/	0.81	0.99	/	0.89	0.93	0.91	0.85	0.96
	silent	< C-L >								0.65	0.90
	Initiate	Start speaking		0.62	0.93		0.43	0.99		0.52	0.96
	speaking	< S-S >		0.02	0.93		0.43	0.99		0.32	0.90
User intention detection	Stop	Start listening	0.91	0.75	0.99	0.86	0.53	0.93	0.89	0.62	0.96
	speaking	< S-L >								0.02	0.90
	Keep	Continue speaking		0.94	1.00		0.96	1.00		0.95	1.00
	speaking	< C-S >		0.94	1.00		0.90	1.00		0.93	1.00

Table 3: Testing using real recorded data with manually labeled user query completion.

		AcousticVAD			+ SemanticVAD					
Threshold	GT\Est	<   C-L   >	< S-S >		<   C-L   >	< S-S >	Recall	Precision	F1	Accuracy
300 ms	< C-L >	0	532	⇒	495	37	0.930	0.887	0.908	0.935
	< S-S >	0	1008		63	945	0.937	0.962	0.949	
500 ms	< C-L >	0	324		307	17	0.949	0.892	0.920	0.962
	< S-S >	0	1078		37	1041	0.966	0.984	0.975	
800 ms	< C-L >	0	228		218	10	0.956	0.861	0.905	0.966
	< S-S >	0	1091		35	1056	0.968	0.991	0.979	0.900
1800 ms	< C-L >	0	132		128	4	0.970	0.800	0.880	0.971
	< S-S >	0	1108		32	1076	0.971	0.996	0.983	0.9/1

we included the corresponding uninterrupted dialogues alongside full-duplex samples in the training set.

For fine-tuning, we used a small version of Hunyuan [28], 0.5B-dense-8k model, on the curated dataset. The initial study focuses on Chinese dialogues and the dataset exclusively consists of Chinese full-duplex dialogues. The training follows standard LLM fine-tuning with four added control tokens, running for 1500 steps with a batch size of 128 and a learning rate of 0.001, linearly decayed to 0.0001 for stable optimization.

## 4. Experimental results

#### 4.1. Test sets

As no benchmark datasets exist for full-duplex evaluation, we generate 1,000 test samples for each interaction scenario <sup>1</sup>. Following Sect. 3.3, we reconstruct the topic and speaking style pools, and generate 2,000 multi-round conversations. From these, 1,000 samples per scenario are randomly selected to evaluate the semantic VAD's control token predictions.

#### 4.2. Evaluation results

The evaluation results presented in Table 1 confirm the effectiveness of the proposed semantic VAD in predicting control tokens across all scenarios. Notably, the detection performance for user barge-in (determining whether the system should continue speaking or switch to start listening) is exceptionally high, benefiting from contextual cues during simultaneous speech, making it more stable. In contrast, user state detection (distinguishing between finished and unfinished queries) is slightly lower as it relies solely on semantic completeness, which varies with speaking style and linguistic nuances, introducing ambiguity and making it inherently more challenging for the proposed method.

Table 2 compares our approach with two recently proposed full-duplex SDS methods [18, 23]. As we do not have access

to their model checkpoints or test datasets, we can only present the results reported in their studies directly for reference. Both studies framed user state and barge-in detection as classification tasks based on acoustic and linguistic patterns with limited semantic understanding. In contrast, our semantic VAD leverages LLM-driven semantic comprehension, enabling more context-aware and reliable full-duplex interaction, achieving relatively better performance.

## 4.3. Evaluation on real-recordings

We further tested our system using internal recordings from user interactions with a half-duplex SDS, focusing on user state detection (start-speaking and continue-listening). Users were instructed to include natural hesitations to simulate scenarios with incomplete queries. The recordings were labeled using both VAD techniques and manual annotation for evaluation. Table 3 presents the results. Since acoustic VAD relies solely on acoustic information and determines the completion of the query by comparing the silence duration against a fixed threshold (e.g. 300 ms), it can only predict <|S-S|>—assuming the user has finished speaking. Semantic VAD refinement significantly improved accuracy, achieving over 93.5% across all cases. Further error analysis indicated that some mispredictions were introduced due to ASR errors rather than limitations of the VAD itself.

## 5. Conclusion

This paper introduces a semantic VAD-based dialogue manager (DM) for full-duplex SDS, leveraging a fine-tuned 0.5B LLM to regulate turn-taking through control tokens. The proposed approach effectively distinguishes between intentional and unintentional barge-ins, detects query completion, and reduces computational overhead by selectively invoking the core dialogue engine (CDE). Experimental results demonstrate improved interaction fluidity and intent recognition. Future work will address system delay and robustness while extending this method to support large multimodal models.

<sup>&</sup>lt;sup>1</sup>Test data can be accessed at https://github.com/semanticVAD/testsets

## 6. References

- [1] K. Jokinen and M. McTear, *Spoken dialogue systems*. Morgan & Claypool Publishers, 2009.
- [2] O. Lemon and O. Pietquin, "Machine learning for spoken dialogue systems," in European Conference on Speech Communication and Technologies (Interspeech'07), 2007, pp. 2685–2688.
- [3] J. Edlund, J. Gustafson, M. Heldner, and A. Hjalmarsson, "Towards human-like spoken dialogue systems," *Speech communica*tion, vol. 50, no. 8-9, pp. 630–645, 2008.
- [4] L. Zhou, J. Gao, D. Li, and H.-Y. Shum, "The design and implementation of xiaoice, an empathetic social chatbot," *Computational Linguistics*, vol. 46, no. 1, pp. 53–93, 2020.
- [5] OpenAI, "ChatGPT can now see, hear, and speak," https://openai. com/blog/chatgpt-can-now-see-hear-and-speak., 2023, [Online; accessed 2025-02-11].
- [6] —, "Introducing ChatGPT," https://openai.com/blog/chatgpt# OpenAI, 2023, [Online; accessed 2025-02-11].
- [7] M. U. Hadi, R. Qureshi, A. Shah, M. Irfan, A. Zafar, M. B. Shaikh, N. Akhtar, J. Wu, S. Mirjalili et al., "Large language models: a comprehensive survey of its applications, challenges, limitations, and future prospects," *Authorea Preprints*, 2023.
- [8] M. F. McTear, "Spoken dialogue technology: enabling the conversational user interface," ACM Computing Surveys (CSUR), vol. 34, no. 1, pp. 90–169, 2002.
- [9] Q. Zhou, B. Li, L. Han, and M. Jou, "Talking to a bot or a wall? how chatbots vs. human agents affect anticipated communication quality," *Computers in Human Behavior*, vol. 143, p. 107674, 2023.
- [10] Y. Mou and K. Xu, "The media inequality: Comparing the initial human-human and human-ai social interactions," *Computers in Human Behavior*, vol. 72, pp. 432–440, 2017.
- [11] H. Sacks, E. A. Schegloff, and G. Jefferson, "A simplest systematics for the organization of turn-taking for conversation," *language*, vol. 50, no. 4, pp. 696–735, 1974.
- [12] G. Skantze, "Turn-taking in conversational systems and humanrobot interaction: a review," *Computer Speech & Language*, vol. 67, p. 101178, 2021.
- [13] D. H. Zimmermann and C. West, "Sex roles, interruptions and silences in conversation," in *Amsterdam Studies in the Theory and History of Linguistic Science Series 4*. John Benjamins BV, 1996, pp. 211–236.
- [14] C. Liu, J. Jiang, C. Xiong, Y. Yang, and J. Ye, "Towards building an intelligent chatbot for customer service: Learning to respond at the appropriate time," in *Proceedings of the 26th ACM SIGKDD* international conference on Knowledge Discovery & Data Mining, 2020, pp. 3377–3385.
- [15] M. Marge, C. Espy-Wilson, N. G. Ward, A. Alwan, Y. Artzi, M. Bansal, G. Blankenship, J. Chai, H. Daumé III, D. Dey et al., "Spoken language interaction with robots: Recommendations for future research," Computer Speech & Language, vol. 71, p. 101255, 2022.
- [16] C. Wang, H. Pan, Y. Liu, K. Chen, M. Qiu, W. Zhou, J. Huang, H. Chen, W. Lin, and D. Cai, "Mell: Large-scale extensible user intent classification for dialogue systems with meta lifelong learning," in *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, 2021, pp. 3649–3659.
- [17] D. Shin, "Llm-based natural conversational agent with speech collision detection for early prompt abort," 2024.
- [18] T.-E. Lin, Y. Wu, F. Huang, L. Si, J. Sun, and Y. Li, "Duplex conversation: Towards human-like interaction in spoken dialogue systems," in *Proceedings of the 28th ACM SIGKDD Conference* on Knowledge Discovery and Data Mining, 2022, pp. 3299–3308.
- [19] P. Wang, S. Lu, Y. Tang, S. Yan, W. Xia, and Y. Xiong, "A full-duplex speech dialogue scheme based on large language models," arXiv preprint arXiv:2405.19487, 2024.

- [20] X. Zhang, Y. Chen, S. Hu, X. Han, Z. Xu, Y. Xu, W. Zhao, M. Sun, and Z. Liu, "Beyond the turn-based game: Enabling real-time conversations with duplex models," arXiv preprint arXiv:2406.15718, 2024.
- [21] C. Fu, H. Lin, Z. Long, Y. Shen, M. Zhao, Y. Zhang, S. Dong, X. Wang, D. Yin, L. Ma et al., "Vita: Towards open-source interactive omni multimodal llm," arXiv preprint arXiv:2408.05211, 2024.
- [22] A. Défossez, L. Mazaré, M. Orsini, A. Royer, P. Pérez, H. Jégou, E. Grave, and N. Zeghidour, "Moshi: a speech-text foundation model for real-time dialogue," arXiv preprint arXiv:2410.00037, 2024
- [23] L. Mai and J. Carson-Berndsen, "Real-time textless dialogue generation," arXiv preprint arXiv:2501.04877, 2025.
- [24] H. Zhang and D. Wang, "Deep learning for acoustic echo cancellation in noisy and double-talk scenarios," *Training*, vol. 161, no. 2, p. 322, 2018.
- [25] ——, "Neural cascade architecture for multi-channel acoustic echo suppression," *IEEE/ACM Transactions on Audio, Speech,* and Language Processing, vol. 30, pp. 2326–2336, 2022.
- [26] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray et al., "Training language models to follow instructions with human feedback," Advances in neural information processing systems, vol. 35, pp. 27730–27744, 2022.
- [27] Yuanbao, "Tencent yuanbao chat," https://yuanbao.tencent.com/ chat/, 2025, [Online; accessed 2025-02-11].
- [28] X. Sun, Y. Chen, Y. Huang, R. Xie, J. Zhu, K. Zhang, S. Li, Z. Yang, J. Han, X. Shu et al., "Hunyuan-large: An open-source moe model with 52 billion activated parameters by tencent," arXiv preprint arXiv:2411.02265, 2024.