

Mind-Paced Speaking: A Dual-Brain Approach to Real-Time Reasoning in Spoken Language Models

Anonymous ACL submission

Abstract

Real-time Spoken Language Models (SLMs) struggle to leverage Chain-of-Thought (CoT) reasoning due to the prohibitive latency of generating the entire thought process sequentially. Enabling SLMs to think while speaking, similar to humans, is attracting increasing attention. We present, for the first time, Mind-Paced Speaking (MPS), a brain-inspired framework that enables high-fidelity, real-time reasoning. Similar to how humans utilize distinct brain regions for thinking and responding, we propose a novel dual-brain approach, employing a "Formulation Brain" for high-level reasoning to pace and guide a separate "Articulation Brain" for fluent speech generation. This division of labor eliminates mode-switching, preserving the integrity of the reasoning process. Experiments show that MPS significantly outperforms existing think-while-speaking methods and achieves reasoning performance comparable to models that pre-compute the full CoT before speaking, while drastically reducing latency. Under a zero-latency configuration, the proposed method achieves an accuracy of 92.8% on the mathematical reasoning task Spoken-MQA and attains a score of 82.5 on the speech conversation task URO-Bench. Our work effectively bridges the gap between high-quality reasoning and real-time interaction.

1 Introduction

Speech has emerged as a more natural and fundamental modality for human-computer interaction, leading to growing emphasis on spoken language models (SLMs) (Cui et al., 2024; Wu et al., 2025; Hu et al., 2025; Yu et al., 2024; Défossez et al., 2024). These models facilitate seamless communication by processing and generating audio-based inputs and outputs. A key component enhancing their capability is the integration of thinking, particularly through Chain-of-Thought (CoT) processes and its extensions (Wei et al., 2022; Wang et al.,

2023; Yao et al., 2023; Gao et al., 2023), as implemented in frameworks like Think-Before-Speak (TBS) (Wei et al., 2022; Xie et al., 2025a; Guo et al., 2025). This approach enables models to decompose complex tasks into step-by-step reasoning sequences, thereby improving interpretability and performance in dialogue systems.

However, generating complete CoT sequences often introduces significant latency, which hinders real-time applications. Recent efforts to reduce reasoning latency have garnered significant attention (Chiang et al., 2025; Xie et al., 2025c). These methods explore "think-while-speaking" paradigms, where models interleave thinking and response tokens. The Large Language Model (LLM) continuously switches between think and response modes. It first generates several think tokens, then produces several response tokens based on them. These response tokens are sent to the Text To Speech (TTS) system for speech synthesis. While the speech is synthesizing, the LLM continues to generate more think tokens. However, this interleaving disrupts semantic coherence by forcing the model to frequently switch between thinking and response generation, potentially degrading the performance.

In fact, the human brain provides a biological analogy for efficient parallel processing. Cognitive neuroscience reveals that thinking and speaking involve distinct brain areas (Nersessian, 2002; Hickok and Poeppel, 2007). Speech does not follow a rigid "think-then-speak" sequence or interleaved sequence. Crucially, it exhibits an incremental nature where later parts of a thought are still being processed while the initial parts of the utterance are already being spoken (Indefrey, 2011). Inspired by this, we introduce Mind-Paced Speaking (MPS), a novel architecture for enabling SLMs to "think" and "speak" in a concurrent and integrated manner. The core of MPS is a dual-brain framework that operates analogously to the human cognitive-speech

084	system. One LLM acts as a central "Formulation Brain", continuously generating an internal stream of thought. The other functions as an Articulation Brain, which receives this thought stream in segments and generates the corresponding spoken output. The Formulation Brain does not need to complete a full reasoning chain before the Articulation Brain begins. Instead, the ongoing thinking process actively sets the pace and provides the contextual guidance for the Articulation Brain, allowing it to vocalize fluently even as the underlying thoughts are still being formed and refined by the Formulation Brain. This mind-paced mechanism ensures that the spoken output is not only grounded in a thinking process but also maintains semantic coherence, closely mimicking the natural human process of thinking while speaking. Furthermore, we propose a think-incomplete Supervised Fine-Tuning (SFT) method to enable the Articulation Brain to respond based on incomplete thinking content. The experimental results on benchmarks such as mathematical reasoning, dialogue, and question-answering, prove that compared to methods that answer directly without thinking, or existing methods that think while speaking, the proposed MPS method effectively utilizes the thinking process and continuous semantic context, obtaining more accurate and higher-quality responses. Compared to TBS method, the proposed MPS significantly reduces response latency while maintaining performance.	136
085		137
086		138
087		139
088		140
089		141
090		142
091		143
092		144
093		
094		145
095		
096		146
097		147
098		148
099		149
100		150
101		151
102		152
103		153
104		154
105		155
106		156
107		157
108		158
109		159
110		160
111		161
112		162
113		163
114		164
115	Our main contribution can be summarized as follows:	165
116		166
117	(1) We propose an MPS architecture that enables SLMs to achieve human-like think-while-speaking capabilities. This method significantly reduces the latency of the CoT process while maintaining the semantic coherence of the LLM. Consequently, the LLM leverages the CoT content to deliver superior performance.	167
118		168
119		
120		
121		
122		
123		
124	(2) We develop a think-incomplete SFT to train LLMs to generate responses based on partial thinking processes, thereby enabling them to perform think-while-speaking.	169
125		170
126		171
127		172
128		173
129		174
130		175
131		176
132		177
133		178
134		179
135		180
		181
		182
		183
		184

these studies remain limited to audio-in-text-out Audio LLMs, not SLMs that can engage in dialogue with humans. In (Wu et al., 2025), Step-Audio 2, which takes speech as input and output, using CoT and reinforcement learning to improve the response qualities, is proposed. Step-Audio 2 offers a solution for introducing explicit reasoning into SLMs. Some methods achieve simultaneous thinking and speaking by segmenting CoT content and response content, using the LLM to generate interleaved think tokens and response tokens (Chiang et al., 2025; Xie et al., 2025c). However, this approach differs from the LLM’s original response generation format. The LLM needs to continuously switch between think mode and response mode, which disrupts semantic coherence and affects its performance.

3 Method

This section first outlines the conventional TBS-based SLM. We then present the proposed MPS method. We also introduce the think-incomplete SFT, which is designed to teach LLMs the think-while-speaking capability.

3.1 Think Before Speaking

The architecture of TBS-based SLM is shown in Figure 1. To enhance the reasoning ability of SLM, the TBS paradigm, after receiving user speech \mathbf{X}^{spc} and optional text instructions \mathbf{X}^{txt} , first generates step-by-step CoT tokens $\mathbf{Y}^{\text{cot}} \in R^{T_c}$, and then generates the response tokens $\mathbf{Y}^{\text{res}} \in R^{T_r}$, where T_c and T_r denote the number of CoT tokens and response tokens, respectively. This can be divided into two processes: the thinking process and the speaking process. The thinking process can be written as:

$$P_{\theta_l}(\mathbf{Y}^{\text{cot}} | \langle \mathbf{X}^{\text{spc}}, \mathbf{X}^{\text{txt}} \rangle) = \prod_{t=1}^{T_c} P_{\theta_l}(Y_t^{\text{cot}} | \langle \mathbf{Y}_{1:t-1}^{\text{cot}}, \mathbf{X}^{\text{spc}}, \mathbf{X}^{\text{txt}} \rangle), \quad (1)$$

where θ_l denotes the parameters of the SLM. After that, the LLM generates response tokens for speaking, which can be formulated as:

$$P_{\theta_l}(\mathbf{Y}^{\text{res}} | \langle \mathbf{Y}^{\text{cot}}, \mathbf{X}^{\text{spc}}, \mathbf{X}^{\text{txt}} \rangle) = \prod_{t=1}^{T_r} P_{\theta_l}(Y_t^{\text{res}} | \langle \mathbf{Y}_{1:t-1}^{\text{res}}, \mathbf{Y}_{1:T_c}^{\text{cot}}, \mathbf{X}^{\text{spc}}, \mathbf{X}^{\text{txt}} \rangle). \quad (2)$$

Through this method, the task is decomposed into a step-by-step process. Additionally, by introducing CoT tokens, it enables more Transformer forward operations and thus gives LLM a deeper inference depth (Goyal et al., 2023; Pfau et al., 2024).

3.2 Architecture

In the human brain, speech production is not a monolithic process but the result of two highly specialized and collaborative systems. The first, a network centered around the prefrontal-temporal cortex, is responsible for high-level cognitive functions such as conceptualization, logical reasoning, and content planning. Subsequently, a second system, primarily involving the motor cortex and sub-cortical pathways, translates these abstract thoughts into natural language for articulation, enabling fluent speech. These two systems operate in parallel, with the cognitive system continuously supplying thinking content to the articulatory system, creating a natural flow where the mind paces speech (Nersessian, 2002; Hickok and Poeppel, 2007).

Inspired by this, we abstract this mechanism of separated "formulation" and "articulation" into our model architecture. Instead of relying on a single LLM to handle both thinking and speaking, we propose a dual-brain system composed of two distinct LLMs.

Our proposed framework, illustrated in Figure 2, leverages a dual-LLM architecture consisting of a Formulation Brain LLM and an Articulation Brain LLM. The Formulation Brain LLM is dedicated to user intent understanding and performs deliberate CoT reasoning, with its internal process materialized as "think tokens". Subsequently, the Articulation Brain LLM converts this structured reasoning and the dialogue context into natural language, producing the final "response tokens" for spoken output.

Formulation Brain: The Formulation Brain’s operating mode is identical to that of TBS Audio LLMs but with only the thinking process. After receiving user input \mathbf{X}^{spc} and \mathbf{X}^{txt} , it aims to generate the step-by-step CoT tokens $\mathbf{Y}^{\text{cot}} \in R^{T_c}$. We use the tokens $\langle \text{think} \rangle$ and $\langle / \text{think} \rangle$ to mark the beginning and end of the CoT. This process can be formulated as (1). In the MPS architecture, we do not wait for the Formulation Brain to complete the entire CoT before the Articulation Brain starts speaking. We divide the CoT tokens \mathbf{Y}^{cot} into N segments, denoted as $[\mathcal{S}_1^{\text{cot}}, \mathcal{S}_2^{\text{cot}}, \dots, \mathcal{S}_N^{\text{cot}}]$. Each time the Formulation Brain produces a think seg-

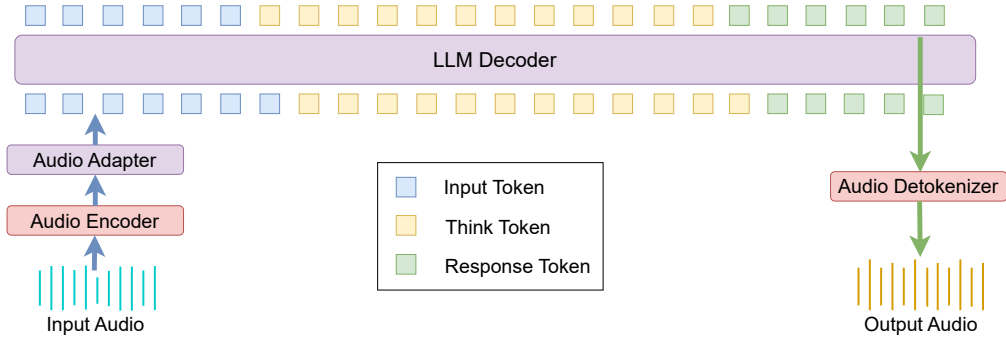


Figure 1: Architecture of the TBS architecture. For the sake of conciseness, we remove the input text, which is optional in SLMs. The TBS SLM first generates the full CoT and then produces response tokens.

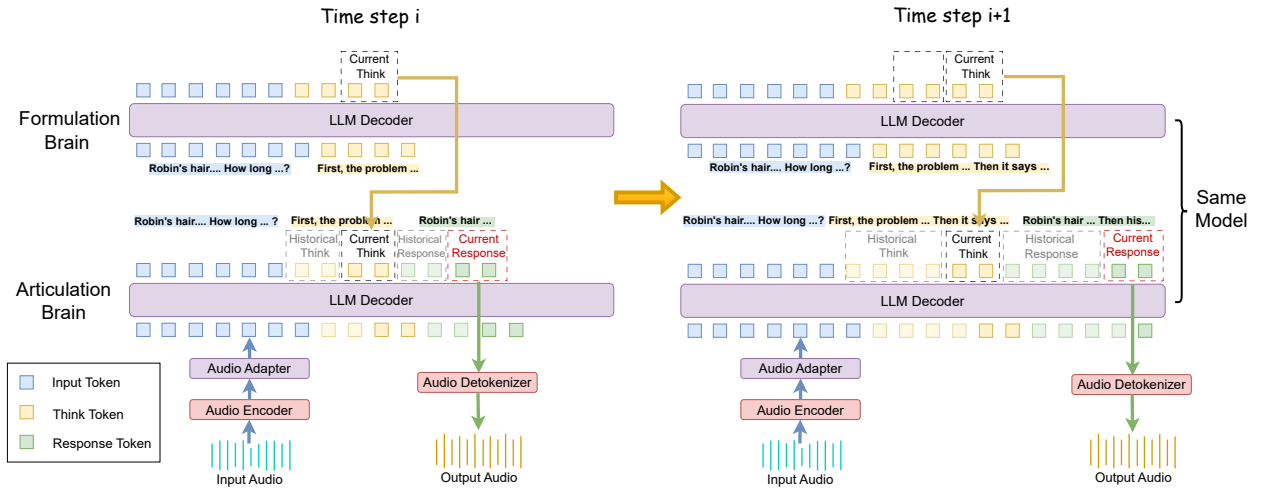


Figure 2: Architecture of the proposed MPS. For the sake of conciseness, we remove the input text, which is optional in SLMs. We demonstrate the process from step i to step $i+1$ when generating think segments and response segments. The Formulation Brain LLM continuously generates the think segments. The newly generated think segment and the response segment from the previous step are both added as the prefix to the Articulation Brain LLM, pacing the Articulation Brain LLM to produce response segment correspondingly.

276 ment S_n^{cot} , we feed the segment to the Articulation
 277 Brain, which then generates a response segment
 278 based on the current think segment and historical
 279 thinking and response contents. After the Formu-
 280 lation Brain LLM finishes CoT segments, it stops
 281 generating response tokens as we do not require
 282 the Formulation Brain to speak.

283 **Articulation Brain:** The Articulation Brain ac-
 284 cepts the same user input as the Formulation
 285 Brain. After obtaining the current think segment S_n^{cot}
 286 from the Formulation Brain, we concatenate it with
 287 the historical think segments $[S_1^{\text{cot}}, S_2^{\text{cot}}, \dots, S_{n-1}^{\text{cot}}]$,
 288 placing $\langle \text{think} \rangle$ and $\langle / \text{think} \rangle$ at the begin-
 289 ning and the end, and then append the histor-
 290 ical response segments, which are defined as
 291 $[S_1^{\text{res}}, S_2^{\text{res}}, \dots, S_{n-1}^{\text{res}}]$. This allows the Articulation
 292 Brain to continue generating the subsequent re-

293 sponse content. After that, we use a streaming TTS
 294 model to synthesize speech in real-time. The Ar-
 295 ticulation Brain’s output is incremental. For every
 296 think segment that the Formulation Brain produces,
 297 the Articulation Brain generates a segment of the
 298 response S_n^{res} . This process can be written as:

$$\begin{aligned}
 P_{\theta_i}(S_n^{\text{res}} | \langle S^{\text{cot}}, X^{\text{spc}}, X^{\text{txt}} \rangle) = \\
 \prod_{n=1}^N P_{\theta_i}(S_n^{\text{res}} | \langle S_{1:n-1}^{\text{res}}, S_{1:n}^{\text{cot}}, X^{\text{spc}}, X^{\text{txt}} \rangle).
 \end{aligned}
 \tag{3}$$

300 When the Formulation Brain just begins its think-
 301 ing, the Articulation Brain can only generate a re-
 302 sponse segment based on a small amount of CoT.
 303 The response segment it generates at this stage may
 304 be of lower quality. As the Formulation Brain’s

305 thinking content increases, the Articulation Brain
306 receives more CoT content, and it subsequently
307 generates responses of increasingly higher quality.

308 Compared with existing think-while-speaking
309 methods that use a single LLM to predict inter-
310 leaved think and response tokens, thereby forcibly
311 interrupting and splitting the originally continuous
312 think and response content (Chiang et al., 2025; Xie
313 et al., 2025c), our method adopts a dual-brain de-
314 sign consisting of the Formulation Brain and the Ar-
315 ticulation Brain. From the perspectives of the For-
316 mulation Brain and the Articulation Brain, both are
317 classic TBS LLMs that, after receiving user input,
318 first generate step-by-step CoT content and then
319 generate response content conditioned on the CoT,
320 thereby greatly ensuring the semantic coherence
321 of the LLM output. By allowing the Formulation
322 Brain to pace the Articulation Brain, our method
323 achieves a human-like think-while-speaking pro-
324 cess.

3.3 Think-incomplete SFT

325 Since the proposed MPS method does not change
326 the input-output patterns of the classic LLM for
327 the individual Formulation Brain and Articulation
328 Brain, the proposed MPS, unlike existing think-
329 while-listening methods (Chiang et al., 2025; Xie
330 et al., 2025c), does not require retraining the
331 LLM. To ensure that the Articulation Brain LLM
332 possesses the ability to accept incomplete think
333 content and produce reasonable output, we intro-
334 duce think-incomplete SFT. In the construction of
335 training data, we randomly retain the content of
336 the first L steps of the step-by-step CoT, delete
337 the subsequent CoT content, then place this in-
338 complete CoT with `<think>` and `</think>` tokens
339 at the beginning and end, concatenate it with the
340 groundtruth response, and use it as the next-token-
341 prediction training objective for the LLM (Brown
342 et al., 2020).
343

344 During the inference stage, we use segments
345 with a fixed number of tokens. We set T_c and T_r to
346 80 and 100 respectively. We use the output format
347 of Step-Audio 2, specifically the *ta4* format, which
348 generates one text token followed by four speech
349 tokens, thus every 100 response tokens contain 20
350 text tokens and 80 speech tokens. The speech syn-
351 thesis frame rate is 12.5 Hz, meaning that 80 tokens
352 correspond to 6.4 seconds of speech. This dura-
353 tion is fully sufficient for the LLM to generate 80
354 think tokens, thereby ensuring the continuous and
355 uninterrupted speech generation of our proposed

method. We also attempt to use the same segment
356 division strategy as in the think-incomplete SFT
357 phase, but we find that it does not bring improve-
358 ment; on the contrary, it introduces uncontrollable
359 latency due to the variable length of each CoT step.
360 We also try using a fixed token count strategy for
361 dividing the CoT during the think-incomplete SFT
362 phase, but it does not yield performance improve-
363 ments either.
364

4 Experiments

4.1 Experimental Settings

365 The LLM backbone used in this paper is Step-
366 Audio 2, and its parameter settings refer to (Wu
367 et al., 2025). The LLMs in Formulation Brain and
368 Articulation brain share the same parameters. To
369 verify the effectiveness of the proposed method on
370 tasks requiring reasoning, we use Spoken-MQA, a
371 mathematical reasoning dataset (Wei et al., 2025).
372 We use accuracy as the evaluation metric. Fur-
373 thermore, to validate the method’s effectiveness on
374 general dialogue tasks, we introduce URO-Bench,
375 which contains several subtasks such as daily di-
376 alogue, emotion recognition, paralinguistic infor-
377 mation, and question-answering (Yan et al., 2025).
378 For question-answering tasks, we use accuracy as
379 the metric. For other tasks, we use GPT-score, gen-
380 erated by GPT-4o-mini and ranging from 0 to 100,
381 to evaluate response quality.
382

383 To accommodate the latency requirements of
384 different application scenarios, we implement two
385 distinct MPS paradigms:
386

- 387 • Think-First, denoted as **MPS-thkfirst**: The
388 Formulation Brain LLM first generates T_c
389 think tokens, after which the Articulation
390 Brain LLM generates T_r response tokens and
391 synthesizes speech. Under this setting, the
392 latency is T_c plus the buffer size of stream-
393 ing TTS, which is significantly lower than the
394 latency required for the TBS structure to gen-
395 erate a complete Chain-of-Thought.
- 396 • Speak-First, denoted as **MPS-spkfirst**: The
397 Articulation Brain LLM first generates T_r re-
398 sponse tokens, while simultaneously, the For-
399 mulation Brain LLM begins generating think
400 tokens. The Formulation Brain LLM com-
401 pletes generating T_c think tokens before the
402 speech synthesized from the T_r response to-
403 kens finishes playing. In this configuration,
404 the latency is solely the buffer size of the

Table 1: Test-set accuracy (%) of different methods on the Spoken-MQA benchmark. The evaluated approaches include: the direct response baseline without a thinking process (MPS-wo), Think-Before-Speaking (MPS-tbs), Think-First (MPS-thkfirst), and Speak-First (MPS-spkfirst). Results of baseline systems are taken from (Xie et al., 2025c) except that results of Step-Audio 2 are reproduced by ourselves.

Method	Arithmetic			Reasoning			Avg
	Short	Long	Avg	Single	Multi	Avg	
Whisper-Qwen2.5-7B-Instruct	-	-	70.0	-	-	72.5	72.2
Whisper-Qwen2.5-Math-7B-Instruct	-	-	77.3	-	-	86.7	85.6
LLaMA-Omni	40.0	11.0	23.5	29.5	10.5	16.2	16.8
Mini-Omni	5.0	2.3	3.5	0.8	1.9	1.6	1.7
Freeze-omni	43.0	14.5	26.8	69.0	19.8	34.4	33.3
GLM-4-Voice	40.0	22.5	30.1	54.4	28.5	36.2	35.3
Qwen2-Audio-7B-Instruct	43.0	31.2	36.3	55.4	22.5	32.3	32.7
Qwen2.5-Omni-7B	83.0	45.1	61.5	85.2	71.5	75.6	73.6
Qwen2.5-Omni-3B	84.0	43.3	60.1	81.5	57.1	64.4	63.6
Mini-Omni-Reasoner	92.9	66.1	77.3	85.9	60.5	68.1	68.6
Step-Audio 2	89.0	52.6	65.9	95.6	90.4	91.9	88.8
MPS-wo <thk< th=""></thk<>	71.0	34.1	47.6	88.0	67.8	73.8	70.6
MPS-tbs	90.0	88.4	89.0	94.4	93.2	93.6	93.0
MPS-thkfirst	89.0	84.9	86.4	95.6	94.6	94.9	93.9
MPS-spkfirst	87.0	71.7	77.3	96.0	94.5	94.9	92.8

streaming TTS, meaning the model can be considered to respond directly with near-zero latency.

Additionally, we compare the proposed method with two approaches that use the same LLM backbone as in this paper: Think-Before-Speaking (MPS-tbs) and direct response without thinking (MPS-wo), to validate the effectiveness of our proposed think-while-speaking methodology. All results reported are from a single run with the random seed fixed to 42.

4.2 Data Construction

We begin with real-world user queries as our seed set. To ensure topical diversity and sufficient scale, we employ GPT-4o (OpenAI, 2024) for transcription and augmentation of these queries. These augmented queries are then used as user prompts to distill dialogue data with native CoT from the DeepSeek-R1 model (Guo et al., 2025). Finally, the speech is synthesized by Step-Audio-TTS-3B (et al., 2025).

However, the raw data generated by DeepSeek-R1, a text-centric model, presents two critical challenges for spoken dialogue

Table 2: The average accuracy of different models with CoT capability on Spoken-MQA, and the extra tokens generated by the model before generating the first response token. The evaluated approaches include: Interleaved Think-While-Speaking (Mini-Omni-Reasoner), Think-Before-Speaking (MPS-tbs), Think-First (MPS-thkfirst), and Speak-First (MPS-spkfirst).

Method	Accuracy	Extra Tokens
Mini-Omni-Reasoner	68.6%	8
MPS-tbs	93.0%	762
MPS-thkfirst	93.9%	80
MPS-spkfirst	92.8%	0

applications: (1) Text-specific stylizations, such as Markdown formatting and emojis, which are incompatible with speech synthesis. (2) The CoT data reflects complete, turn-based reasoning chains, a format unsuitable for training the model to respond from partial thoughts. When the CoT generated by the LLM exhibits some incomplete, its performance is affected.

To address these challenges, we implement a

Table 3: Performance of different methods on the URO-Bench. The evaluated approaches include: the direct response baseline without a thinking process (MPS-wo/thk), Think-Before-Speaking (MPS-tbs), Think-First (MPS-thkfirst), and Speak-First (MPS-spkfirst). Results of baseline systems are taken from (Wu et al., 2025). The results of *Multilingual* of URO-Bench are included in *English*.

Method	Language	Basic				Pro			
		U.	R.	O.	Avg	U.	R.	O.	Avg
GPT-4o Audio	Chinese	89.4	65.5	85.2	78.6	70.6	57.2	70.2	67.1
GPT-Realtime		88.8	72.9	90.8	80.6	72.3	62.6	74.2	70.6
Kimi-Audio		79.3	64.7	79.8	73.6	60.4	59.3	76.2	66.0
Qwen-Omni		59.7	69.7	77.3	69.0	59.0	59.8	58.7	59.1
Step-Audio 2		91.1	75.5	86.1	83.3	74.8	63.2	65.1	68.3
MPS-wo/thk	Chinese	91.6	77.3	87.7	83.4	75.1	74.7	72.9	74.4
MPS-tbs		92.6	82.4	93.8	87.8	75.3	84.2	79.5	79.0
MPS-thkfirst		93.6	84.0	94.8	89.1	75.2	84.2	85.2	80.5
MPS-spkfirst		92.5	82.5	93.1	87.6	77.2	84.8	79.0	79.9
GPT-4o Audio	English	90.2	75.9	90.4	84.5	60.7	64.4	78.5	67.5
GPT-Realtime		87.4	84.1	94.1	88.1	59.7	74.5	76.1	68.9
Kimi-Audio		83.4	42.3	60.4	60.0	50.3	40.6	56.0	49.8
Qwen-Omni		66.3	69.6	76.2	70.6	44.5	63.9	49.4	51.0
Step-Audio 2		92.7	76.5	84.9	83.9	64.9	67.8	66.3	66.1
MPS-wo/thk	English	91.5	68.7	78.8	77.4	73.4	79.2	55.8	65.1
MPS-tbs		92.3	81.5	87.5	86.1	76.4	86.4	87.1	83.3
MPS-thkfirst		94.2	81.4	89.0	87.0	76.5	89.3	89.4	85.0
MPS-spkfirst		94.1	78.5	87.5	85.2	76.0	89.7	69.9	74.8

fine-grained data processing pipeline:

- **Compatibility Processing:** We discard samples containing Markdown formatting or multi-item lists that cannot be naturally rendered in speech. For samples containing emojis, we employ Qwen-72B-Instruct (et al., 2024) to remove these elements while preserving the plain text content.
- **CoT Pruning:** To train the model to respond stably with only partial CoT, we augment the data by randomly deleting some reasoning paragraphs. This operation is performed in a way that generally preserves the overall logic of the CoT. To maintain the stylistic distribution of the original DeepSeek CoT, we neither delete individual sentences within a paragraph nor use an LLM to rewrite the content of the remaining parts. This ensures that the preserved paragraphs are stylistically and distributionally consistent with the source model.

4.3 Results

4.3.1 Evaluation on Reasoning Tasks

Table 1 shows the computational accuracy on Spoken-MQA. It can be seen that the proposed MPS-thkfirst method exceeds the MPS-wo/thk method and all baseline methods, including the think-while-speaking Mini-Omni-Reasoner, in all evaluation tasks. The results proves that the proposed method effectively utilizes the thinking process, achieving more intelligent response. Compared to Mini-Omni-Reasoner, the MPS method maintains semantic coherence and achieves better performance.

Besides, compared to MPS-tbs, the MPS-thkfirst method demonstrates comparable performance, being slightly weaker in arithmetic computation tasks but superior in reasoning tasks. One possible explanation is that reasoning tasks require more textual analysis. The MPS-thkfirst method, by using each think segment to pace the generation of a corre-

478	sponding response segment, implicitly achieves	530
479	semantic alignment, enabling the model to better	531
480	utilize contextual information for response genera-	
481	tion.	
482	The MPS-spkfirst method experiences some per-	
483	formance degradation because it initially outputs a	
484	response segment without utilizing any think seg-	
485	ments. This impact is particularly pronounced in	
486	tasks involving direct arithmetic computation. For	
487	reasoning tasks, experimental observations indicate	
488	that the initial phase of the LLM’s CoT content	
489	primarily involves analyzing the semantic informa-	
490	tion of the question, often rewriting the question’s	
491	content. Consequently, this initial portion of the	
492	CoT content has a limited effect on the final re-	
493	sponse. As a result, MPS-spkfirst is minimally	
494	affected in reasoning tasks, and its performance re-	
495	mains nearly identical to that of MPS-thkfirst. Ex-	
496	perimental results on Spoken-MQA demonstrate	
497	that the proposed method significantly leverages	
498	CoT to achieve more intelligent responses. Fur-	
499	thermore, compared to TBS methods, our think-	
500	while-speaking approach achieves comparable per-	
501	formance, with significantly lower CoT latency as	
502	analysed in Section 4.1.	
503	To demonstrate the latency differences between	
504	the proposed method and baseline approaches, we	
505	select four models from Table 1 that include a think-	
506	ing process: Mini-Omni-Reasoner, MPS-tbs, MPS-	
507	thkfirst, and MPS-spkfirst. We calculate the num-	
508	ber of extra tokens generated by the model from the	
509	end of the user’s question to the generation of the	
510	first response token on Spoken-MQA. The results	
511	are shown in Table 2. It can be observed that com-	
512	pared to MPS-tbs, MPS-thkfirst achieves higher	
513	accuracy while exhibiting significantly lower re-	
514	sponse latency. Although the accuracy of MPS-	
515	spkfirst is slightly lower than that of MPS-tbs	
516	and MPS-thkfirst, its response is without latency.	
517	Furthermore, compared to Mini-Omni-Reasoner,	
518	which uses interleaved think and response tokens to	
519	achieve think-while-speaking, the proposed MPS	
520	methods achieve higher accuracy. Notably, the	
521	MPS-spkfirst attains this superior accuracy with	
522	zero latency. This indicates that MPS-spkfirst can	
523	play a more critical role in real-time dialogue sce-	
524	narios with low-latency requirements. To further	
525	demonstrate how the Formulation Brain drives the	
526	Articulation Brain to realize the process of thinking	
527	while speaking, we present and describe an exam-	
528	ple of MPS-spkfirst on the SpokenMQA dataset in	
529	Appendix B.	
	4.3.2 Evaluation on Speech-to-speech	
	Conversation	
	Table 3 shows the results of different methods	532
	on URO-Bench. It can be observed that MPS-	533
	thkfirst achieves higher performance than MPS-	534
	tbs on nearly all tasks and on average, under lower	535
	response latency. This may also be related to the	536
	implicit semantic alignment performed by MPS-	537
	thkfirst. Due to generating an initial response	538
	segment without prior thinking, MPS-spkfirst per-	539
	forms slightly worse than MPS-thkfirst, but still sig-	540
	nificantly outperforms the direct response method	541
	MPS-wo/thk. It can also be observed that this ini-	542
	tial response segment without thinking, primarily	543
	affects the performance of MPS-spkfirst on Oral	544
	Conversation within the <i>Pro</i> dataset. In contrast,	545
	the performance on Understanding, Reasoning,	546
	and the subtasks in all the <i>Basic</i> dataset is less	547
	affected. Nevertheless, MPS-spkfirst features low-	548
	est response latency as analysed in Section 4.1, and	549
	its response performance remains close to or even	550
	better than that of MPS-tbs in some tasks, making	551
	it more suitable for scenarios requiring faster feed-	552
	back. The experimental results demonstrate that	553
	the proposed MPS method maintains high perfor-	554
	mance on dialogue tasks, achieving performance	555
	comparable to TBS models while operating at sig-	556
	nificantly lower latency.	557
	5 Conclusion	558
	This paper proposes the MPS method, which en-	559
	ables SLMs to possess the ability to think while	560
	speaking. Inspired by the human thinking and re-	561
	sponse mechanism, we use a Formulation Brain	562
	LLM to continuously generate think segments, pac-	563
	ing the Articulation Brain LLM to utilize historical	564
	and current think segments, as well as historical	565
	responses, to generate current response segment.	566
	Experimental results on mathematical reasoning	567
	and speech conversation tasks show that the pro-	568
	posed method significantly outperforms direct re-	569
	sponse methods and existing think-while-speaking	570
	methods. It achieves performance comparable or	571
	even better than methods that complete thinking	572
	before responding, while greatly reducing response	573
	latency. The proposed method breaks through the	574
	limitations of existing interleaved thinking and	575
	response-based think-while-speaking methods and	576
	provides an effective reference for researching real-	577
	time dialogue consistent with human thinking and	578
	response mechanisms.	579

580 Limitations

581 Despite the breakthroughs achieved in real-time
582 reasoning for SLMs and in enabling LLMs to
583 "think while listening" in a human-like manner,
584 our work presents several limitations:

- 585 • Resource Consumption: Our method requires
586 two LLMs for inference. While the utilization
587 of the prefix technique ensures that the total
588 number of inference steps remains same to
589 a single LLM, it doubles the memory usage
590 for two LLMs. Additionally, the dual-model
591 architecture introduces communication over-
592 head between the components.
- 593 • Hyperparameter Optimization: The number
594 of tokens assigned to the think and response
595 segments is determined based on empirical
596 results and theoretical time overhead calcula-
597 tions. However, the theoretical optimum for
598 these configurations remains to be fully ex-
599 plored.
- 600 • Generalizability: Our investigation is con-
601 ducted exclusively using the Step-Audio LLM
602 backbone. The generalizability of our ap-
603 proach across other LLM architectures re-
604 mains to be verified through further experi-
605 ments.

606 References

607 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
608 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
609 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
610 Askell, and 1 others. 2020. Language models are
611 few-shot learners. *Advances in neural information*
612 *processing systems*, 33:1877–1901.

613 Cheng-Han Chiang, Xiaofei Wang, Linjie Li, Chung-
614 Ching Lin, Kevin Lin, Shujie Liu, Zhendong Wang,
615 Zhengyuan Yang, Hung-yi Lee, and Lijuan Wang.
616 2025. Stitch: Simultaneous thinking and talking
617 with chunked reasoning for spoken language models.
618 *arXiv preprint arXiv:2507.15375*.

619 Wenqian Cui, Dianzhi Yu, Xiaoqi Jiao, Ziqiao Meng,
620 Guangyan Zhang, Qichao Wang, Yiwen Guo, and Ir-
621 win King. 2024. Recent advances in speech language
622 models: A survey. *arXiv preprint arXiv:2410.03751*.

623 Alexandre Défossez, Laurent Mazaré, Manu Orsini,
624 Amélie Royer, Patrick Pérez, Hervé Jégou, Edouard
625 Grave, and Neil Zeghidour. 2024. Moshi: a speech-
626 text foundation model for real-time dialogue. *arXiv*
627 *preprint arXiv:2410.00037*.

Ailin Huang et al. 2025. *Step-audio: Unified under- 628*
standing and generation in intelligent speech interac- 629
tion. *Preprint*, arXiv:2502.11946. 630

An Yang et al. 2024. *Qwen2 technical report*. *Preprint*, 631
arXiv:2407.10671. 632

Qingkai Fang, Yan Zhou, Shoutao Guo, Shaolei Zhang, 633
and Yang Feng. 2025. Llama-omni2: Llm-based real- 634
time spoken chatbot with autoregressive streaming 635
speech synthesis. *arXiv preprint arXiv:2505.02625*. 636

Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, 637
Pengfei Liu, Yiming Yang, Jamie Callan, and Gra- 638
ham Neubig. 2023. Pal: Program-aided language 639
models. In *International Conference on Machine 640*
Learning, pages 10764–10799. PMLR. 641

Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Kr- 642
ishna Menon, Sanjiv Kumar, and Vaishnavh Nagara- 643
jan. 2023. Think before you speak: Training lan- 644
guage models with pause tokens. *arXiv preprint 645*
arXiv:2310.02226. 646

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao 647
Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi- 648
rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. 649
Deepseek-r1: Incentivizing reasoning capability in 650
llms via reinforcement learning. *arXiv preprint 651*
arXiv:2501.12948. 652

Gregory Hickok and David Poeppel. 2007. The cortical 653
organization of speech processing. *Nature reviews 654*
neuroscience, 8(5):393–402. 655

Ke Hu, Ehsan Hosseini-Asl, Chen Chen, Edresson 656
Casanova, Subhankar Ghosh, Piotr Żelasko, Zhehuai 657
Chen, Jason Li, Jagadeesh Balam, and Boris Gins- 658
burg. 2025. Efficient and direct duplex modeling for 659
speech-to-speech language model. *arXiv preprint 660*
arXiv:2505.15670. 661

Peter Indefrey. 2011. The spatial and temporal sig- 662
natures of word production components: a critical 663
update. *Frontiers in psychology*, 2:255. 664

Gang Li, Jizhong Liu, Heinrich Dinkel, Yadong Niu, 665
Junbo Zhang, and Jian Luan. 2025. Reinforcement 666
learning outperforms supervised fine-tuning: A case 667
study on audio question answering. *arXiv preprint 668*
arXiv:2503.11197. 669

Nancy J Nersessian. 2002. *The cognitive basis of model- 670*
-based reasoning in science. na. 671

OpenAI. 2024. *Gpt-4o system card*. *Preprint*, 672
arXiv:2410.21276. 673

Jacob Pfau, William Merrill, and Samuel R Bowman. 674
2024. Let’s think dot by dot: Hidden computa- 675
tion in transformer language models. *arXiv preprint 676*
arXiv:2404.15758. 677

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, 678
Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan 679
Zhang, YK Li, Yang Wu, and 1 others. 2024. 680

681	Deepseekmath: Pushing the limits of mathematical reasoning in open language models. <i>arXiv preprint arXiv:2402.03300</i> .	736
682		737
683		738
684	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models . In <i>The Eleventh International Conference on Learning Representations</i> .	739
685		740
686		741
687		742
688		743
689		744
690	Chengwei Wei, Bin Wang, Jung-jae Kim, and Nancy F Chen. 2025. Towards spoken mathematical reasoning: Benchmarking speech-based models over multi-faceted math problems. <i>arXiv preprint arXiv:2505.15000</i> .	745
691		746
692		747
693		748
694		749
695	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	750
696		751
697		752
698		753
699		754
700		755
701	Cheng Wen, Tingwei Guo, Shuaijiang Zhao, Wei Zou, and Xiangang Li. 2025. Sari: Structured audio reasoning via curriculum-guided reinforcement learning. <i>arXiv preprint arXiv:2504.15900</i> .	756
702		757
703		758
704		759
705	Boyong Wu, Chao Yan, Chen Hu, Cheng Yi, Chengli Feng, Fei Tian, Feiyu Shen, Gang Yu, Haoyang Zhang, Jingbei Li, and 1 others. 2025. Step-audio 2 technical report. <i>arXiv preprint arXiv:2507.16632</i> .	760
706		761
707		762
708		763
709	Jingran Xie, Shun Lei, Yue Yu, Yang Xiang, Hui Wang, Xixin Wu, and Zhiyong Wu. 2025a. Leveraging chain of thought towards empathetic spoken dialogue without corresponding question-answering data. In <i>ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 1–5. IEEE.	764
710		765
711		766
712		767
713		768
714		769
715		770
716	Zhifei Xie, Mingbao Lin, Zihang Liu, Pengcheng Wu, Shuicheng Yan, and Chunyan Miao. 2025b. Audio-reasoner: Improving reasoning capability in large audio language models. <i>arXiv preprint arXiv:2503.02318</i> .	771
717		772
718		773
719		774
720		775
721	Zhifei Xie, Ziyang Ma, Zihang Liu, Kaiyu Pang, Hongyu Li, Jialin Zhang, Yue Liao, Deheng Ye, Chunyan Miao, and Shuicheng Yan. 2025c. Mini-omni-reasoner: Token-level thinking-in-speaking in large speech models. <i>arXiv preprint arXiv:2508.15827</i> .	776
722		777
723		778
724		779
725		780
726		781
727	Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan, Kai Dang, and 1 others. 2025. Qwen2. 5-omni technical report. <i>arXiv preprint arXiv:2503.20215</i> .	782
728		783
729		784
730		785
731	Ruiqi Yan, Xiquan Li, Wenxi Chen, Zhikang Niu, Chen Yang, Ziyang Ma, Kai Yu, and Xie Chen. 2025. Uro-bench: A comprehensive benchmark for end-to-end spoken dialogue models. <i>arXiv preprint arXiv:2502.17810</i> .	786
732		787
733		788
734		789
735		790
	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. <i>Advances in neural information processing systems</i> , 36:11809–11822.	791
		792
	Wenyi Yu, Siyin Wang, Xiaoyu Yang, Xianzhao Chen, Xiaohai Tian, Jun Zhang, Guangzhi Sun, Lu Lu, Yuxuan Wang, and Chao Zhang. 2024. Salmonn-omni: A codec-free llm for full-duplex speech understanding and generation. <i>arXiv preprint arXiv:2411.18138</i> .	793
		794
	Aohan Zeng, Zhengxiao Du, Mingdao Liu, Kedong Wang, Shengmin Jiang, Lei Zhao, Yuxiao Dong, and Jie Tang. 2024. Glm-4-voice: Towards intelligent and human-like end-to-end spoken chatbot. <i>arXiv preprint arXiv:2412.02612</i> .	795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900

Question:

"Robin's hair was 16 inches long. He cut off 11 inches. Then it again grew by 12 inches. How long is his hair now?"

Groundtruth:

"17"

Response:

Step1:

Formulation Brain:

""

Articulation Brain:

"Robin's hair started at 16 inches. He cut off 11"

Step1:

Formulation Brain:

"First, the problem states that Robin's hair was initially 16 inches long. He cut off 11 inches. So, after cutting, the length would be 16 inches minus 11 inches. That gives 5 inches. Then"

Articulation Brain:

"Robin's hair started at 16 inches. He cut off 11 inches, so 16 minus 11 is 5 inches. Then, his "

Step2:

Formulation Brain:

"First, the problem states that Robin's hair was initially 16 inches long. He cut off 11 inches. So, after cutting, the length would be 16 inches minus 11 inches. That gives 5 inches. Then, it says his hair grew by 12 inches. So, starting from the 5 inches after cutting, adding 12 inches would be $5 + 12 = 17$ inches. The problem says "

Articulation Brain:

"Robin's hair started at 16 inches. He cut off 11 inches, so 16 minus 11 is 5 inches. Then, his hair grew by 12 inches, so 5 plus 12 equals "

Step3:

Formulation Brain:

"First, the problem states that Robin's hair was initially 16 inches long. He cut off 11 inches. So, after cutting, the length would be 16 inches minus 11 inches. That gives 5 inches. Then, it says his hair grew by 12 inches. So, starting from the 5 inches after cutting, adding 12 inches would be $5 + 12 = 17$ inches. The problem says \"he cut off 11 inches,\" which clearly indicates cutting, so subtracting is correct. So, the final answer should be 17 inches."

Articulation Brain:

"Robin's hair started at 16 inches. He cut off 11 inches, so 16 minus 11 is 5 inches. Then, his hair grew by 12 inches, so 5 plus 12 equals 17 inches. Therefore, Robin's hair is now 17 inches long. "

Figure 3: An example of the output of MPS-spkgfirst on the Spoken-MQA dataset. The Articulation Brain first generates a response segment. Simultaneously, Formulation Brain continuously generates new think segments, and each newly generated think segment is prefixed to the Articulation Brain, pacing it to generate new response segment.