

Efficient and Direct Duplex Modeling for Speech-to-Speech Language Model

Anonymous preprint

Abstract

Spoken dialogue is the most intuitive form of human-computer interaction, yet current speech language models often remain constrained to turn-based exchanges, lacking real-time adaptability such as user barge-in. We propose a novel duplex speech to speech (S2S) architecture featuring continuous user inputs and codec agent outputs with channel fusion that directly models simultaneous user and agent streams. Using a pretrained streaming encoder for user input enables the first duplex S2S model without requiring speech pretrain. Separate architectures for agent and user modeling facilitate codec fine-tuning for better agent voices and halve the bitrate (0.6 kbps) compared to previous works. Experimental results show that the proposed model outperforms previous duplex models in reasoning, turn-taking, and barge-in abilities. The model requires significantly less speech data, as speech pretrain is skipped, which markedly simplifies the process of building a duplex S2S model from any LLMs. Finally, it is the first openly available duplex S2S model with training and inference code to foster reproducibility.

Index Terms: duplex, speech-to-speech, conversation, barge-in

1. Introduction

Large language models (LLMs) [1–4] have made significant strides in natural language processing, sparking interest in multimodal models that extend beyond text. Speech, as a natural interface for human-computer interaction, is a key part of this trend. Recent studies suggest adapting LLMs to process speech prompts for various speech-to-text (STT) tasks [2, 4–9].

While traditional systems often respond with text, speech outputs are more intuitive for human-computer interaction. Cascaded spoken dialogue systems, like AudioGPT [10], use text as an intermediate representation, involving sequential modules such as ASR, LLM, and TTS. However, these systems face drawbacks like high latency, lack of interactive behaviors, and loss of paralinguistics. To address these issues, research has shifted towards end-to-end speech-to-speech (S2S) modeling.

Previous S2S models focus on half-duplex, turn-based interactions. For instance, SpeechGPT [11], initialized from LLaMA, undergoes sequential fine-tuning on speech-only data and multimodal instruction sets to handle spoken question-answer (QA) tasks. Similarly, USDM [12] extends Mistral’s pretraining with interleaved speech-text data for enhanced multimodal understanding. GLM-4-voice [13] efficiently tokenizes speech using one codebook and large-scale speech-text pretraining for downstream tasks like ASR, TTS, and SQA.

Several pioneering or concurrent full-duplex S2S models have been recently proposed [14–17]. However, these systems face increased complexity in model, data, and computation, which hinders their widespread research and adoption.

The introduction of additional submodules for turn-taking between user and agent increases system complexity and reduces the end-to-end nature of the models. Moreover, the extensive speech-text pretraining required on top of the LLM backbone is resource-intensive and limits scalability to any LLMs. Finally, using codecs to model user and agent interactions simultaneously necessitates a delicate balance between speech perception and generation, presenting another significant challenge.

To tackle the above problems, we propose a novel duplex S2S system with the following contributions: 1) A novel duplex S2S architecture featuring continuous user inputs and codec agent outputs with channel fusion that directly models simultaneous text and speech of both the user and agent. 2) We demonstrate several key advantages over existing duplex models: The use of a pretrained encoder as input enables the first duplex S2S model **without speech pretraining** requirement; As the agent and user are modeled by the codec and the pretrained encoder separately, this facilitates **codec fine-tuning** toward better agent voices. 3) We propose a set of systematic metrics to evaluate conversational behaviors such as turn-taking and barge-in. Finally, it is the first **open duplex S2S model** with both training and inference code publicly available to foster reproducibility.

2. Related Work

Interest in full-duplex S2S models has grown in the past year. Key challenges here include handling simultaneous user and agent streams and enabling turn-taking. Systems like [16, 18, 19] model single-channel interactions but use external signals, such as stopping commands [19] or submodules [16], to decide when to respond. Models like SyncLLM [20] and OmniFlatten [17] achieve full-duplex conversation by employing time chunking methods, embedding time information into LLMs for synchronization. This interleaving processing allows the model to handle user inputs like barge-in with low latency.

Our duplex S2S model is trained without speech-text pretraining, unlike [14]. In multi-turn conversation, we align text and speech at the turn level, which simplifies data preparation compared to word-level alignment. Compared to [15, 18, 19], our model predicts text and speech simultaneously without requiring an explicit TTS component. Our speech codec model uses parallel codebooks (see details in Sec. 3.2) and enables speech generation with minimal latency. Our design further enables codec fine-tuning for improved agent voices while halving the required bitrate of previous works (0.6 kbps).

3. Model Architecture

To achieve duplex behavior, our S2S model takes two input streams simultaneously: user speech stream, and agent speech and text stream. As shown in Fig. 1, the user speech is first

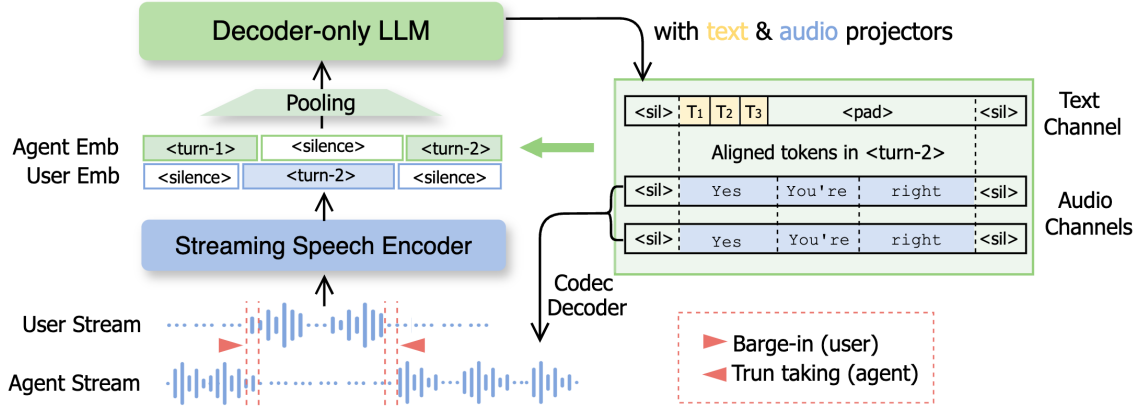


Figure 1: The proposed duplex S2S model **without** requiring speech-text pretraining. Our model includes a streaming speech encoder, a personalized codec, and an LLM. The model is trained to predict both text and audio channels in parallel with turn-level alignments.

96 encoded to generate continuous embeddings by the speech encoder using an 80-ms frame rate. We use a 100M streaming
 97 encoder from a CTC model [21]. We initialize the backbone LLM using the TinyLlama-1.1B-chat model [22]. A
 98 modality adapter is used between the speech encoder and the text LLM. To obtain the agent embeddings in training, we use
 99 a codec model [23] to generate 12.5 Hz speech codes for the agent speech. LLM vocabulary is extended to include extra
 100 tokens from speech codec with zero initialization. The two inputs are time-aligned and summed as the input to the text LLM
 101 (similar to [24]). Both our speech encoder and text LLM are causal and thus streaming. In training, we fine-tune both the
 102 speech encoder and the backbone LLM. Text and speech loss are weighted differently in training (see Sec. 5.1). Our model is
 103 trained by multi-channel next token prediction similar to [1].

111 3.1. Simultaneous Agent Text and Speech Prediction

112 As shown in Fig. 1, we encode speech using 4 codebooks at a rate of 12.5 frames per second [23], and text targets are
 113 tokenized into a separate channel. We align the text and speech tokens at the turn level based on their start time. We prepend
 114 separate <BOS> tokens for text and speech at the beginning of the turn and append <EOS> at the end of the turn. The gap
 115 between text and speech tokens are padded by text pad ID. We also tried word-level alignment between text and speech as in [14]
 116 and did not find improvement. Empirically, we find that the model tends to learn agent text first. Therefore, we introduce a
 117 small delay (i.e., one token) to the speech channels for improved
 118 speech quality without introducing significant latency.

124 3.2. Personalization-friendly Speech Tokenization

125 We employ a partially causal neural audio codec to transform
 126 raw speech signals into streaming tokenized representations.
 127 Given an audio signal \mathbf{a} , the codec generates a two-dimensional
 128 acoustic matrix, $\mathbf{C}_{T \times N} = \text{CodecModel}(\mathbf{a})$, where T denotes
 129 the downsampled sequence length, and N represents the number
 130 of codebooks per timestep. Each element in $\mathbf{C}_{T \times N}$ is an
 131 m -bit discrete code. We adopt the state-of-the-art NanoCodec
 132 [23], which achieves reasonable-quality audio compression at
 133 0.6 kbps with a frame rate of 12.5 frames per second, employ-
 134 ing $N = 4$ independent codebooks. The codec leverages Finite
 135 Scalar Quantization (FSQ) [25], ensuring independence among
 136 codebooks. This independence removes the need for additional
 137 models or delay mechanisms, allowing all N codebooks to be

Table 1: Synthetic training data with multi-turn and barge-in.

Task	Dataset	#Hours	Speech	Multi-turn	Barge-in
Spoken QA	ASR-QA	20k	Mix	Augment	×
	MS MARCO	0.2k	TTS	Augment	×
	Alpaca	0.2k	TTS	Augment	×
	Internal SFT	3k	TTS	Real	✓
Convers- ation	UltraChat	3k	TTS	Augment	✓
	Topic	0.3k	TTS	Augment	✓



Figure 2: Duplex training data format. Our duplex data consists of separate user and agent streams including turn taking and barge-in behavior. Here, the user barges in at the second turn.

138 predicted in parallel at each timestep, thereby enabling fully
 139 parallel modeling with low latency.

140 Our duplex design allows us to personalize the pretrained
 141 codec for agent voices to further enhance audio quality. This
 142 is enabled by modeling the agent and user separately with the
 143 speech codec and a pretrained causal speech encoder. In the
 144 experimental section, we will evaluate the benefits of speech
 145 and reasoning quality resulting from codec personalization.

146 4. Duplex Data for Training

147 Table 1 summarizes our training data which can be categorized
 148 into spoken QA and multi-turn conversations.

149 4.1. Single-turn synthetic and real spoken QA

150 Our most basic training data structure consists of a single-turn
 151 spoken QA between the user and agent. We use a multi-speaker
 152 TTS model [26] to synthesize the context, questions and answers
 153 from *MS MARCO* [27] and *Alpaca* [28]. To mitigate
 154 overfitting to synthetic data, we follow [29] to create additional
 155 synthetic QA pairs using the Mixtral-8x22B LLM from English
 156 ASR-labeled data (8k public¹ and 12k in-house). This data
 157 is then synthesized using the same TTS, denoted as *ASR-QA*.

¹A subset from the NeMo ASR set in [21].

158 The resultant user speech contains both TTS and real data. An
159 evaluation set used in Sec. 5.2 is created from the public data
160 portion. We use a fixed speaker to generate agent speech and
161 randomly select speakers for user speech.

162 We create duplex training data from the aforementioned
163 user-agent QA pairs. First, we split a pair of utterances into
164 two streams, corresponding to the user and agent portions sep-
165 arately, and then insert silence into the agent stream when the
166 user speaks, and vice versa. This gives us two streams of speech
167 (shown as the first turn in Fig.2). This duplex structure enables
168 the model to listen and speak simultaneously at any time. To
169 prevent the agent from barging in unexpectedly, we insert a
170 0.64s silence between user and agent before the agent speaks.

171 4.2. Augment with Multi-turn and Barge-in

172 In order for the model to learn the ability for multi-turn con-
173 versation, we also create duplex data that includes two or more
174 turns of conversation between the user and agent (e.g., Fig. 2).
175 First, we synthesize 3k hours of duplex data from a text-based
176 multi-turn *Internal SFT* dataset to form multi-turn spoken QA.
177 To ensure a more conversational flow, we limit each turn of the
178 text SFT data, which is typically very long, to under 25 seconds.
179 Second, we augment the single-turn data from Sec. 4.1 by ran-
180 domly concatenating two QA pairs from the same dataset. The
181 multi-turn data topics focus on role-playing, daily topics, scien-
182 tific topics, etc. Moreover, when creating multi-turn data, we
183 allow the user to barge in by cutting off the agent speech. After
184 the cutoff, we keep a small duration (0.64 s) of agent speech to
185 account for barge-in latency, and pad the rest of the agent turn
186 with silence. As we show in later results, this straightforward
187 approach enables the model to learn barge-in behavior.

188 4.3. Conversational data

189 To enhance the model’s conversation ability on daily topics,
190 we create *Topic* and *UltraChat* datasets (totaling 3.3k hours as
191 shown in Table 1). For both datasets, we first generate 4-turn
192 text-based conversations and then synthesize them using a TTS
193 model [26]. For *Topic*, we randomly choose a topic between
194 user and agent and prompt the Meta-Llama-3.1-70B-Instruct
195 model [4] to generate a conversation. The topics are randomly
196 chosen from the everyday-conversation dataset [30], which cov-
197 ers 63 everyday and science topics. To generate concise replies
198 for efficient training, we restrict the words of each turn to be
199 30 words in the prompt. The generated conversations are then
200 synthesized into speech and prepared to the duplex data format.
201 For *UltraChat*, we randomly sample a chat conversation from
202 the *UltraChat* dataset [31] to use as contextual information in
203 the prompt to generate a 4-turn conversation similar to *Topic*.

204 5. Experiment Details

205 5.1. Training Details

206 We implement the model with PyTorch using the NeMo Toolkit
207 [32], and the model is trained on 32 A100 (80G) GPUs with a
208 batch duration of 1000 sec per GPU. The speech encoder is ini-
209 tialized from a 100M streaming pretrained encoder with 80ms
210 right context [21], and the LLM is initialized from the 1.1B
211 TinyLlama [22]. We use a 32k SentencePiece tokenizer for text,
212 and a personalized 0.6 kbps NanoCodec [23] for speech by de-
213 fault. Ablations for personalization are presented in Sec. 6.3.
214 The speech codes have 4 channels, with a vocabulary size of
215 4,037 for each channel. Text and speech channel training loss

are weighted by 3 and 1 respectively. We use FusedAdam, and
an inverse Square Root Annealing learning rate (LR) schedule
for optimization. The LR schedule starts with an initial learning
rate of $3e-4$ with a warm-up of 2500 steps. Gradient clipping is
applied at the threshold of 1.0 to stabilize training.

221 5.2. Evaluation Data and Metrics

222 Our evaluation data consists of: 1) *multi-turn* conversations: *Ul-*
223 *traChat*, *Roleplay* (part of *Internal SFT*), and *Topic*, and 2) *spo-*
224 *ken QA reasoning*: *ASR-QA* and *Alpaca*. We select one shard
225 for each dataset in Sec. 4, which is unseen during training, for
226 this evaluation. To evaluate model performance on a more chal-
227 lenging scenario where the user frequently interrupts the agent,
228 we create an evaluation set called *Impatient*. When creating *Im-*
229 *patient*, we halve the silence time between the current and the
230 next user turn (from the original duration in the *ASR-QA* set) to
231 increase the chance of the agent being interrupted by the user.
232 By doing this, the interruption cases for our model and Moshi
233 (more details in Sec. 6.1) in the *Impatient* dataset are as high as
234 95.4% and 96.7%, respectively.

235 In terms of evaluation metrics, we evaluate the reason-
236 ing ability of our model using GPT scores generated by
237 `gpt-4o-mini-2024-07-18` ranging from 0 to 10 based on the
238 hypotheses and references of all the agent turns. The reason-
239 ing quality is evaluated using the aforementioned *multi-turn*
240 and *spoken QA reasoning* datasets. The hypotheses of agent
241 turns are produced by transcribing the generated speech using
242 the ASR model `nvidia/parakeet-tdt_ctc-110m`.

243 We evaluate turn-taking ability and speech generation qual-
244 ity using the *UltraChat* and *Impatient* datasets. We use two
245 types of metrics to measure the turn-taking ability: barge-in per-
246 formance and 1st response latency (see Table 2). For barge-in
247 performance, we introduce the following metrics: 1) Barge-in
248 latency: The time delay between the user’s speech onset and
249 the agent stopping its response; 2) Success rate: The percent-
250 age of cases where the agent successfully stops speaking within
251 1.5 seconds after the user interruption; and 3) False alarm rate:
252 The frequency at which the agent incorrectly barges in while
253 the user speaks. Additionally, if the user stops speaking within
254 0.1s, the event is not counted as a false alarm, as we found that
255 Moshi tends to proactively respond. The 1st response latency is
256 defined as the time taken by the agent to respond to the 1st user
257 turn. To evaluate the speech quality, we compute the UTMOS
258 [33] using the generated agent speech after removing silence.

259 6. Results and Comparison

260 6.1. Conversation and Speech Generation Quality

261 We first evaluate the turn-taking and speech generation quality
262 of our model in Table 2. Compared to Moshi, our model has
263 significantly higher barge-in success rate (94.5% v.s. 55.1%),
264 the same false alarm rates, and lower barge-in latency (0.69s
265 v.s. 0.81s). We observe that, in multi-turn conversations, Moshi
266 often initiates dialogue more proactively, leading to user barge-
267 in failures for both *UltraChat* and *Impatient*.

268 We cannot directly compare our 1st response latency with
269 Moshi’s as Moshi almost always responds before the user fin-
270 ishes talking and thus does not fit for this metric. We also note
271 that our 1st response latency is affected by our data generation,
272 as we always add a 0.64-second silence after the user turns to
273 ensure no unexpected agent barge-in. Further reducing this de-
274 lay is our future work. Lastly, we report UTMOS and our model
275 generates better quality speech than Moshi by up to 0.4.

Table 2: Comparison of turn-taking and speech generation quality.

Dataset	Model	Barge-in Performance			1st Response	UT
		Success \uparrow	False Alarms \downarrow	Latency (s) \downarrow	Latency (s) \downarrow	MOS \uparrow
UltraChat	Ours	83.0%	0.0%	0.52	0.72	4.3
	Moshi	56.0%	0.0%	0.63	n/a	3.9
Impatient	Ours	94.5%	0.0%	0.69	0.92	4.0
	Moshi	55.1%	0.0%	0.81	n/a	3.8

Table 3: Reasoning quality of multi-turn conversation and spoken QA. GT+LLM denotes an optimal cascaded system which feeds every ground-truth user turn to the LLM.

GPT Score	Multiturn Conversation			Spoken QA	
	UltraChat	Roleplay	Topic	ASR-QA	Alpaca
Ours	3.5	4.6	6.1	7.8	2.9
Moshi	3.4	1.7	2.8	1.9	1.7
GT+LLM	<u>6.4</u>	<u>4.9</u>	5.5	5.8	<u>5.0</u>

Table 4: Evaluation of audio reconstruction and the resultant S2S quality across different codecs.

Codec	Bitrate kbps	Audio Reconstruction			S2S
		MOS \uparrow	CER \downarrow	SECS \uparrow	ASR-BLEU \uparrow
Mimi[14]	1.1	4.16	3.00	0.65	n/a
Nano[23]	1.2	4.67	1.44	0.77	18.1
Nano[23]	0.6	4.54	3.55	0.57	16.2
+personalized	0.6	4.75	1.36	0.94	18.7

6.2. Reasoning Quality

In Table 3, we compare the reasoning ability of our model to Moshi [14] and an *optimal cascaded system* formed by feeding every ground-truth user turn text to LLM (i.e., GT+LLM in Table 3). The backbone of our model, TinyLlama, is used as the LLM. We report the aforementioned GPT scores on two types of test sets: *multi-turn conversation* and *spoken QA*. Compared to Moshi, our model shows better scores on all datasets despite the fact that our model uses much less data and smaller backbone. Compared to the optimal cascaded system, our model shows competitive results, better on two and worse on three sets. The slightly worse performance of end-to-end versus cascaded is not new and has been shown by other research [2, 11, 14, 29]. Future works include i) a more fair comparison with full pipeline (VAD, streaming ASR and TTS, LLM), and ii) improving the reasoning of duplex S2S models.

6.3. Speech Codec Personalization

We personalize the codec to our agent voice by fine-tuning the codec on 21k ground-truth utterances from the target speaker. The model is evaluated on 228 test samples that are not seen during training. Perceptual quality is assessed using estimated Mean Opinion Scores (MOS) with Torchaudio-Squim [34]. Intelligibility is measured by computing the Character Error Rate (CER), comparing transcriptions from the Massively Multilingual Speech (MMS) model [35] for both ground-truth and reconstructed audio. Speaker similarity is evaluated using the Speaker Encoder Cosine Similarity (SECS) [36], computed with the state-of-the-art ECAPA2 speaker encoder [37].

Table 4 presents the evaluation results for the 1.1 kbps Mimi Codec [14], 1.2 kbps, and 0.6 kbps versions of NanoCodec [23], and the proposed personalized version of 0.6 kbps NanoCodec. Personalization significantly enhances the performance of the

0.6 kbps NanoCodec. Notably, despite operating at nearly half the bitrate, our personalized codec outperforms both Mimi and NanoCodec at 1.2 kbps across all audio reconstruction metrics on the target speaker.

As an ablation study, we further train our duplex S2S models with different codecs (last three rows in Table 4). For simplicity, we report ASR-BLEU, which is calculated based on the reference agent texts and ASR transcripts of generated agent speech. Results on *ASR-QA* in Table 4 indicate that personalization enhances duplex modeling as well, leading to improved perceptual quality and higher BLEU scores.

6.4. Listening Duplex Conversation Examples

We include representative listening examples in an anonymous demo page². Specifically, the following capabilities of our duplex S2S model on *unseen* data are highlighted:

Robustness with frequent interruption. In the example of Fig. 3 and the webpage, the user interrupts the agent three times in 15 seconds, and leaves limited time for the agent to respond. Despite these challenges, the agent still demonstrates robust conversational behavior in handling frequent barge-in.



Figure 3: Multi-turn conversation with frequent barge-in.

Unseen reasoning problem. Beyond leveraging learned knowledge to generate responses, the agent also demonstrates the ability to utilize contextual information, effectively summarizing the main topic of each conversation in Fig. 4 and webpage that was unseen during training.

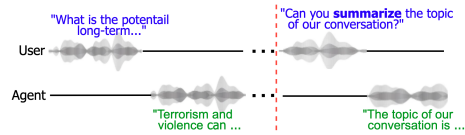


Figure 4: Spoken QA example on an unseen topic.

7. Conclusion

We introduced a novel duplex S2S architecture that models simultaneous user and agent streams without requiring speech pretraining. Our data-efficient approach maintains end-to-end modeling of conversation reasoning and behaviors. Experimental results show competitive performance in reasoning, barge-in, and turn-taking. Our open-sourced training and inference code will also be a valuable resource for future research.

²<https://anonymous598e.github.io/INTERSPEECH2025-DEMO/>

8. References

- 341
342 [1] T. B. Brown, "Language models are few-shot learners," *arXiv preprint arXiv:2005.14165*, 2020.
343
- 344 [2] G. Team, P. Georgiev, V. I. Lei, R. Burnell, L. Bai, A. Gulati,
345 G. Tanzer, D. Vincent, Z. Pan, S. Wang *et al.*, "Gemini 1.5: Un-
346 locking multimodal understanding across millions of tokens of
347 context," *arXiv preprint arXiv:2403.05530*, 2024.
- 348 [3] J. Achiam, S. Adler, S. Agarwal *et al.*, "Gpt-4 technical report,"
349 *arXiv preprint arXiv:2303.08774*, 2023.
- 350 [4] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Let-
351 man, A. Mathur, A. Schelten, A. Yang, A. Fan *et al.*, "The llama
352 3 herd of models," *arXiv preprint arXiv:2407.21783*, 2024.
- 353 [5] Y. Fathullah, C. Wu, E. Lakomkin *et al.*, "Prompting large lan-
354 guage models with speech recognition abilities," in *ICASSP*.
355 IEEE, 2024, pp. 13 351–13 355.
- 356 [6] Y. Chu, J. Xu, X. Zhou, Q. Yang, S. Zhang, Z. Yan, C. Zhou, and
357 J. Zhou, "Qwen-audio: Advancing universal audio understand-
358 ing via unified large-scale audio-language models," *arXiv preprint*
359 *arXiv:2311.07919*, 2023.
- 360 [7] Z. Chen, H. Huang, A. Andrusenko *et al.*, "Salm: Speech-
361 augmented language model with in-context learning for speech
362 recognition and translation," in *ICASSP*. IEEE, 2024, pp. 13 521–
363 13 525.
- 364 [8] Z. Kong, A. Goel, R. Badlani *et al.*, "Audio flamingo: A novel
365 audio language model with few-shot learning and dialogue abili-
366 ties," *arXiv preprint arXiv:2402.01831*, 2024.
- 367 [9] K. Hu, Z. Chen, C.-H. H. Yang, P. Želasko, O. Hrinchuk,
368 V. Lavrukhin, J. Balam, and B. Ginsburg, "Chain-of-
369 thought prompting for speech translation," *arXiv preprint*
370 *arXiv:2409.11538*, 2024.
- 371 [10] R. Huang, M. Li, D. Yang, J. Shi, X. Chang, Z. Ye, Y. Wu,
372 Z. Hong, J. Huang, J. Liu *et al.*, "Audiogpt: Understanding and
373 generating speech, music, sound, and talking head," in *Proceed-*
374 *ings of the AAAI Conference on Artificial Intelligence*, vol. 38,
375 no. 21, 2024, pp. 23 802–23 804.
- 376 [11] D. Zhang, S. Li, X. Zhang, J. Zhan, P. Wang, Y. Zhou, and
377 X. Qiu, "Speechgpt: Empowering large language models with
378 intrinsic cross-modal conversational abilities," *arXiv preprint*
379 *arXiv:2305.11000*, 2023.
- 380 [12] H. Kim, S. Seo, K. Jeong, O. Kwon, J. Kim, J. Lee, E. Song,
381 M. Oh, S. Yoon, and K. M. Yoo, "Unified speech-text pretraining
382 for spoken dialog modeling," *arXiv preprint arXiv:2402.05706*,
383 2024.
- 384 [13] A. Zeng, Z. Du, M. Liu, K. Wang, S. Jiang, L. Zhao, Y. Dong, and
385 J. Tang, "Glm-4-voice: Towards intelligent and human-like end-
386 to-end spoken chatbot," *arXiv preprint arXiv:2412.02612*, 2024.
- 387 [14] A. Défossez, L. Mazaré, M. Orsini, A. Royer, P. Pérez, H. Jégou,
388 E. Grave, and N. Zeghidour, "Moshi: a speech-text foundation
389 model for real-time dialogue," *arXiv preprint arXiv:2410.00037*,
390 2024.
- 391 [15] W. Yu, S. Wang, X. Yang, X. Chen, X. Tian, J. Zhang, G. Sun,
392 L. Lu, Y. Wang, and C. Zhang, "Salmonn-omni: A codec-free
393 llm for full-duplex speech understanding and generation," *arXiv*
394 *preprint arXiv:2411.18138*, 2024.
- 395 [16] Q. Chen, Y. Chen, Y. Chen, M. Chen, Y. Chen, C. Deng, Z. Du,
396 R. Gao, C. Gao, Z. Gao *et al.*, "Minmo: A multimodal large
397 language model for seamless voice interaction," *arXiv preprint*
398 *arXiv:2501.06282*, 2025.
- 399 [17] Q. Zhang, L. Cheng, C. Deng, Q. Chen, W. Wang, S. Zheng,
400 J. Liu, H. Yu, C. Tan, Z. Du *et al.*, "Omniflatten: An end-to-
401 end gpt model for seamless voice conversation," *arXiv preprint*
402 *arXiv:2410.17799*, 2024.
- 403 [18] X. Wang, Y. Li, C. Fu, Y. Shen, L. Xie, K. Li, X. Sun, and L. Ma,
404 "Freeze-omni: A smart and low latency speech-to-speech dia-
405 logue model with frozen llm," *arXiv preprint arXiv:2411.00774*,
406 2024.
- [19] Z. Xie and C. Wu, "Mini-omni2: Towards open-source gpt- 407
408 4o with vision, speech and duplex capabilities," *arXiv preprint*
409 *arXiv:2410.11190*, 2024.
- [20] B. Veluri, B. N. Peloquin, B. Yu, H. Gong, and S. Gollakota, "Be- 410
411 yond turn-based interfaces: Synchronous llms as full-duplex dia-
412 logue agents," *arXiv preprint arXiv:2409.15594*, 2024.
- [21] NVIDIA, "STT En FastConformer Hybrid Transducer- 413
414 CTC Large Streaming 80ms," 2023, version
1.20.0, Released June 22, 2023. [Online]. Avail- 415
416 able: [https://catalog.ngc.nvidia.com/orgs/nvidia/teams/nemo/](https://catalog.ngc.nvidia.com/orgs/nvidia/teams/nemo/models/stt.en.fastconformer.hybrid.large.streaming.80ms)
417 [models/stt.en.fastconformer.hybrid.large.streaming.80ms](https://catalog.ngc.nvidia.com/orgs/nvidia/teams/nemo/models/stt.en.fastconformer.hybrid.large.streaming.80ms)
- [22] P. Zhang, G. Zeng, T. Wang, and W. Lu, "Tinyllama: An open- 418
419 source small language model," 2024.
- [23] Anonymous, "Nanocodec: Towards high-quality ultra fast speech 420
421 llm inference," *Preprint (Under Review)*, 2025.
- [24] Z. Ma, Y. Song, C. Du, J. Cong, Z. Chen, Y. Wang, Y. Wang, 422
423 and X. Chen, "Language model can listen while speaking," *arXiv*
424 *preprint arXiv:2408.02622*, 2024.
- [25] F. Mentzer, D. Minnen, E. Agustsson, and M. Tschannen, "Fi- 425
426 nite scalar quantization: Vq-vae made simple," *arXiv preprint*
427 *arXiv:2309.15505*, 2023.
- [26] S. Hussain, P. Neekhara, X. Yang, E. Casanova, S. Ghosh, M. T. 428
429 Desta, R. Fejgin, R. Valle, and J. Li, "Koel-tts: Enhancing llm
430 based speech generation with preference alignment and classifier
431 free guidance," *arXiv preprint arXiv:2502.05236*, 2025.
- [27] P. Bajaj, D. Campos, N. Craswell, L. Deng, J. Gao, X. Liu, R. Ma- 432
433 jumdar, A. McNamara, B. Mitra, T. Nguyen *et al.*, "A human gen-
434 erated machine reading comprehension dataset," *arXiv preprint*
435 *arXiv:1611.09268*, 2018.
- [28] R. Taori, I. Gulrajani, T. Zhang, Y. Dubois, X. Li, C. Guestrin, 436
437 P. Liang, and T. B. Hashimoto, "Stanford alpaca: An instruction-
438 following llama model," 2023.
- [29] V. Noroozi, Z. Chen, S. Majumdar, S. Huang, J. Balam, and 439
440 B. Ginsburg, "Instruction data generation and unsuper-
441 vised adaptation for speech language models," *arXiv preprint*
442 *arXiv:2406.12946*, 2024.
- [30] H. Face, "Everyday conversations for llms," [https://huggingface.](https://huggingface.co/datasets/HuggingFaceTB/everyday-conversations-llama3.1-2k) 443
444 [co/datasets/HuggingFaceTB/everyday-conversations-llama3.](https://huggingface.co/datasets/HuggingFaceTB/everyday-conversations-llama3.1-2k)
445 [1-2k](https://huggingface.co/datasets/HuggingFaceTB/everyday-conversations-llama3.1-2k), 2024.
- [31] N. Ding, Y. Chen, B. Xu, Y. Qin, Z. Zheng, S. Hu, Z. Liu, 446
447 M. Sun, and B. Zhou, "Enhancing chat language models by
448 scaling high-quality instructional conversations," *arXiv preprint*
449 *arXiv:2305.14233*, 2023.
- [32] O. Kuchaiev, J. Li, H. Nguyen *et al.*, "Nemo: a toolkit for 450
451 building ai applications using neural modules," *arXiv preprint*
452 *arXiv:1909.09577*, 2019.
- [33] T. Saeki, D. Xin, W. Nakata, T. Koriyama, S. Takamichi, and 453
454 H. Saruwatari, "Utmos: Utokyo-sarulab system for voicemos
455 challenge 2022," *arXiv preprint arXiv:2204.02152*, 2022.
- [34] A. Kumar, K. Tan, Z. Ni, P. Manocha, X. Zhang, E. Henderson, 456
457 and B. Xu, "Torchaudio-squim: Reference-less speech quality
458 and intelligibility measures in torchaudio," in *ICASSP 2023-2023*
459 *IEEE International Conference on Acoustics, Speech and Signal*
460 *Processing (ICASSP)*. IEEE, 2023, pp. 1–5.
- [35] V. Pratap, A. Tjandra, B. Shi, P. Tomasello, A. Babu, S. Kundu, 461
462 A. Elkahky, Z. Ni, A. Vyas, M. Fazel-Zarandi *et al.*, "Scaling
463 speech technology to 1,000+ languages," *Journal of Machine*
464 *Learning Research*, vol. 25, no. 97, pp. 1–52, 2024.
- [36] E. Casanova, J. Weber, C. D. Shulby, A. C. Junior, E. Gölge, and 465
466 M. A. Ponti, "Yourtts: Towards zero-shot multi-speaker tts and
467 zero-shot voice conversion for everyone," in *International Con-*
468 *ference on Machine Learning*. PMLR, 2022, pp. 2709–2720.
- [37] J. Thienpondt and K. Demuynck, "Ecapa2: A hybrid neural net- 469
470 work architecture and training strategy for robust speaker embed-
471 dings," *arXiv preprint arXiv:2401.08342*, 2024.