

AV-Dialog: Spoken Dialogue Models with Audio-Visual Input

Anonymous ACL submission

Abstract

Dialogue models falter in noisy, multi-speaker environments, often producing irrelevant responses and awkward turn-taking. We present AV-Dialog, the first multimodal dialog framework that uses both audio and visual cues to track the target speaker, predict turn-taking, and generate coherent responses. By combining acoustic tokenization with multi-task, multi-stage training on monadic, synthetic, and real audio-visual dialogue datasets, AV-Dialog achieves robust streaming transcription, semantically grounded turn-boundary detection and accurate responses, resulting in a natural conversational flow. Experiments show that AV-Dialog outperforms audio-only models under interference, reducing transcription errors, improving turn-taking prediction, and enhancing human-rated dialogue quality. These results highlight the power of seeing as well as hearing for speaker-aware interaction, paving the way for spoken dialogue agents that perform robustly in real-world, noisy environments.

1 Introduction

Dialogue models are moving closer to natural, human-like interaction (Défossez et al., 2024; Veluri et al., 2024), but real-world deployment remains challenging. Real environments are complex with background noise, overlapping talk, and interfering speakers. This setting is known as the “cocktail party problem”: the difficulty of attending to a target speaker amid simultaneous talkers and noise. Current models rely solely on speech inputs, making them brittle in these settings; often losing track of the target speaker, producing irrelevant responses, and breaking natural turn-taking.

We argue that overcoming this limitation requires looking as well as listening. Humans address the “cocktail party problem” by combining auditory and visual cues, using lip movements and gaze to focus on the speaker and learn turn-taking cues (Mcdermott, 2009; Best et al., 2023).

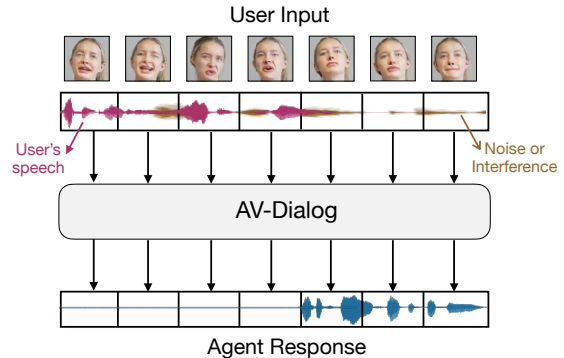


Figure 1: AV-Dialog understands audio-visual input from the target user (purple waveform), accurately detects the appropriate time to take a turn in the conversation, and outputs responses (blue waveform), even in the presence of interfering speakers (brown waveform).

Inspired by this, we present AV-Dialog, a novel audio-visual framework for dialogue modeling. Designing such a framework requires meeting three key challenges: First, the model must continuously process audio and video in a streaming manner, isolating the target speaker even when background noise or louder interfering speakers are present. Second, it must detect turn-taking cues and respond appropriately, maintaining conversational flow despite overlapping or interfering speech. Third, the system must produce coherent responses to the intended speaker without being misled by distractors or environmental noise.

Our paper presents the first spoken dialog models with audio-visual input that address the above challenges. We make the following contributions:

- **Multimodal dialogue modeling.** We start with a pre-trained large language model (LLAMA3-8B) (Dubey et al., 2024) and train it to process audio and video input in a streaming manner with 40ms chunks. The model learns to extract text tokens for the target speaker under interference and predict turn-change tokens for natural conversational timing. We explore two architectures: *dual* and *unified*. In the dual architecture, the multimodal model outputs transcriptions and turn-taking to-

069 kens that trigger a second LLM-based text back- 120
070 bone for high-quality response generation, either 121
071 via in-context learning or instruction-tuning on di- 122
072 alogue data for greater naturalness. The unified 123
073 architecture instead uses a single model to perform 124
074 AV understanding, turn prediction, and response
075 generation jointly. Our results show that explicit
076 turn-change supervision is not only essential for
077 dual-model setups but also improves the generation
078 quality of the unified model.

079 • **Acoustic tokenization for noisy, multi-speaker**
080 **settings.** Unlike prior dialogue models (Défossez 126
081 et al., 2024; Veluri et al., 2024) that rely on seman- 127
082 tic tokenizers (e.g., HuBERT) trained on single- 128
083 speaker speech, we use general-purpose acoustic 129
084 tokens, Descript Audio Codec (DAC) (Kumar et al., 130
085 2023) for multimodal dialogue modeling. Because 131
086 acoustic tokens preserve both semantic and raw 132
087 acoustic information, they enable inherent speaker 133
088 differentiation based on voice characteristics. Thus, 134
089 we can better address the “cocktail party prob- 135
090 lem,” maintaining robustness to noise and inter- 136
091 fering speakers across a range of Signal-to-Noise 137
092 Ratios (SNRs), indicating the noise level of the 138
093 input speech. In ablation studies, replacing seman- 139
094 tic with acoustic tokens both reduces word error 140
095 rate for streaming AVSR from 67% to 31.7% under 141
096 strong multi-speaker interference as well as enables 142
097 more timely responses. 143

098 • **Multi-task, multi-stage training recipe.** Open 144
099 audio-visual dialogue datasets are much smaller 145
100 than text-based chat corpora, making robust train- 146
101 ing challenging. We address this with a multi-task, 147
102 two-stage training strategy: the first stage trains 148
103 the base LLaMA model with text prediction, ASR, 149
104 AVSR and audio captioning tasks to strengthen 150
105 audio-visual understanding and align with original 151
106 text embeddings. The second stage fine-tunes the 152
107 model on real audio-only and audio-visual conversa- 153
108 tional datasets to learn natural turn-taking and 154
109 conversation context. We further improve robust- 155
110 ness with synthetic mixture augmentation, simulat- 156
111 ing noisy, multi-speaker environments. This task- 157
112 oriented approach enables AV-Dialog to acquire 158
113 complementary skills from each dataset, enhancing 159
114 transcription accuracy, turn-taking prediction, and 160
115 dialogue quality under challenging conditions. 161

116 We compare AV-Dialog with Moshi-7B (Défos- 162
117 sez et al., 2024), a state-of-the-art spoken dialogue 163
118 model. Results show that adding the visual modal- 164
119 ity boosts turn-taking prediction accuracy from 165

54% to 79% in the presence of interfering speakers. 120
Human evaluation (N=18) further demonstrates a 121
+1.75-point MOS improvement in dialogue natu- 122
ralness and a +1.99-point MOS gain in response 123
relevance and helpfulness. 124

2 Related work 125

Audio-visual speech recognition. A related task 126
is Audio-Visual Speech Recognition (AVSR) (Rou- 127
ditchenko et al., 2024; Hong et al., 2023), where 128
models like AV-HuBERT (Shi et al., 2022) learn 129
speech representations from synchronized audio 130
and video. Recent work combines pre-trained au- 131
dio (Radford et al., 2023) and video (Shi et al., 132
2022) with language models to improve word er- 133
ror rates (Cappellazzo et al., 2025a). While AVSR 134
systems excel at speech recognition, they are not de- 135
signed for generative dialogue, turn-taking, or full- 136
duplex interaction. In contrast, AV-Dialog extends 137
audio-visual fusion beyond recognition to enable 138
grounded conversational agents. Moreover, most 139
AVSR models operate offline with full-recording 140
access (Rouditchenko et al., 2024; Cappellazzo 141
et al., 2025a), whereas our system performs stream- 142
ing inference and incorporates AVSR as a multi- 143
task objective. We therefore compare AV-Dialog 144
with state-of-the-art streaming AVSR models such 145
as Auto-AVSR (Ma et al., 2023) in §4.1. 146

Dialog models. Recent work on dialog models 147
generate spoken responses from a given prompt. 148
Notably, SpeechGPT (Zhang et al., 2023) is fine- 149
tuned on speech-only data and multimodal instruc- 150
tions for spoken question answering. Multimodal 151
models like SpiritLM (Nguyen et al., 2024) accept 152
speech or text as prompts and generate responses in 153
either modality, while prior non-open source mod- 154
els (Park et al., 2024) handle audio-visual inputs 155
but require explicit prompting. Unlike AV-Dialog, 156
these systems do not model turn-taking and thus 157
do not know when to respond, a key component of 158
human-like dialog interaction. (Liao et al., 2025) 159
leverage audio and visual cues for turn prediction, 160
but their approach is non-streaming, does not sup- 161
port full-duplex interaction, and requires clean text 162
transcripts as input. 163

Recent full-duplex dialogue models like 164
dGSLM (Lakhota et al., 2021), Moshi (Défossez 165
et al., 2024), and SyncLLM (Veluri et al., 2024) 166
generate responses concurrently with user input by 167
predicting intent or turn-endings without explicit 168
prompts. However, relying solely on text and 169

semantic speech tokens limits their ability to track the target speaker in noisy, multi-speaker settings. In contrast, AV-Dialog integrates visual cues and general-purpose acoustic tokens for robust speaker tracking and dialogue generation under challenging signal-to-noise (SNR) conditions.

3 AV-Dialog Models

In human conversation, we process rich acoustic and visual cues to understand and respond via speech. Prevalent dialogue models, however, emphasize the modality the agent must generate (speech & language) while only considering the semantic representations of speech and/or ignore visual cues; making them brittle to noise and interference. In contrast, combining audio and visual context enables accurate turn-boundary detection and timely responses. To this end, we develop a dialogue framework built on audio-visual understanding that infers the user’s complete intent, both what is said and when they intend to yield the floor.

As shown in Fig. 2B, AV-Dialog’s dual-model architecture comprises two components: an audio-visual dialogue understanding module and a text backbone. The latter is implemented with instruction-tuned LLAMA3-8B (Dubey et al., 2024), though any text LLM or API can be used.

3.1 Audio-Visual Dialogue Understanding

We propose an AV dialogue understanding model, fine-tuned on a base text-LLM, with two essential capabilities for voice interfaces: recognizing user speech and detecting intent to yield the conversation floor. It must also operate as a streaming model, so we can update the response LLM’s KV-cache as the user speaks.

Our model processes multi-stream inputs: at each timestep n , a continuous visual stream V_n from a visual encoder and 16 audio streams $\mathbf{A}_n = [A_{n,1}, \dots, A_{n,16}]$ from an audio tokenizer. It outputs two streams: (1) a text stream U_n representing the user’s input, and (2) a turn event stream T_n from the AV understanding module, predicting when the agent should take the conversation floor.

Both streams are synchronized and operate on 40 ms chunks. Both embeddings are projected to the transformer’s model dimension via separate linear layers and summed with the previous timestep’s text embedding to produce the final embedding $e_n = \mathcal{L}_A(\sum_{i=1}^{16} \mathcal{E}(A_{n,i})) + \mathcal{L}_V(V_n) + \mathcal{E}(U_{n-1}) + \mathcal{E}(T_{n-1})$. Here, $\mathcal{E}(\cdot)$ denotes the embedding layer,

$\mathcal{L}_V(\cdot)$ the visual projection layer, and $\mathcal{L}_A(\cdot)$ the audio projection layer.

At the AV-Dialog model output z_n , two linear heads estimate the distributions of U_n and T_n , conditioned on all preceding sub-sequences.

$$\sigma(L_U(z_n)) \approx \mathbb{P}[U_n | \mathbf{A}_{\leq n}, V_{\leq n}, U_{<n}, T_{<n}],$$

$$\sigma(L_T(z_n)) \approx \mathbb{P}[T_n | \mathbf{A}_{\leq n}, V_{\leq n}, U_{<n}, T_{<n}]$$

$\sigma(\cdot)$ is the softmax operation, $L_U(\cdot)$ the linear header for the text stream, and $L_T(\cdot)$ the linear header for the turn event stream.

3.1.1 Audio-Visual Encoding

Most prior turn-taking and spoken dialogue models rely on speaker-invariant semantic speech representations (Nguyen et al., 2022; Veluri et al., 2024; Défossez et al., 2024). While effective in clean settings, they struggle in real-world, “cocktail-party” environments. Models like HuBERT (Hsu et al., 2021), though robust to uncorrelated background noise, can amplify spurious speech interference due to their speaker invariance.

We instead leverage general audio representations, enabling inherent speaker differentiation based on voice characteristics. Our AV-Dialog model uses the high-fidelity Descript Audio Codec (DAC) (Kumar et al., 2023) tokenizer, where each 40 ms chunk is encoded into 16 DAC codebooks.

We incorporate visual cues of the speaker because they (1) enable robust target speech identification in noisy environments, (2) enhance speech perception and understanding (Shi et al., 2022; Cappellazzo et al., 2025b), and (3) provide crucial signals for estimating turn boundaries. Specifically, we use the dlib library (dlib) to detect face regions in first-person video and extract continuous lip-centric visual representations via a pre-trained AV-HuBERT model (Shi et al., 2022).

3.1.2 Output Streams

The AV dialogue understanding module outputs two token streams: i) time-aligned transcription of user’s speech U_n , and ii) turn-taking event labels T_n , using the L_U and L_T heads, respectively.

If a user’s word begins at t_{start} , the model predicts its tokens from timestep $\lceil t_{start}/25 \rceil + d$, where d is a small delay providing a reasonable context for recognizing the word. When no word is uttered, it outputs a silence token $\langle \text{EMP} \rangle$.

For turn boundaries, we adopt Pairwise-TurnGPT’s (Leishman et al., 2024) turn-taking event taxonomy: (1) Normal turn, agent speaks

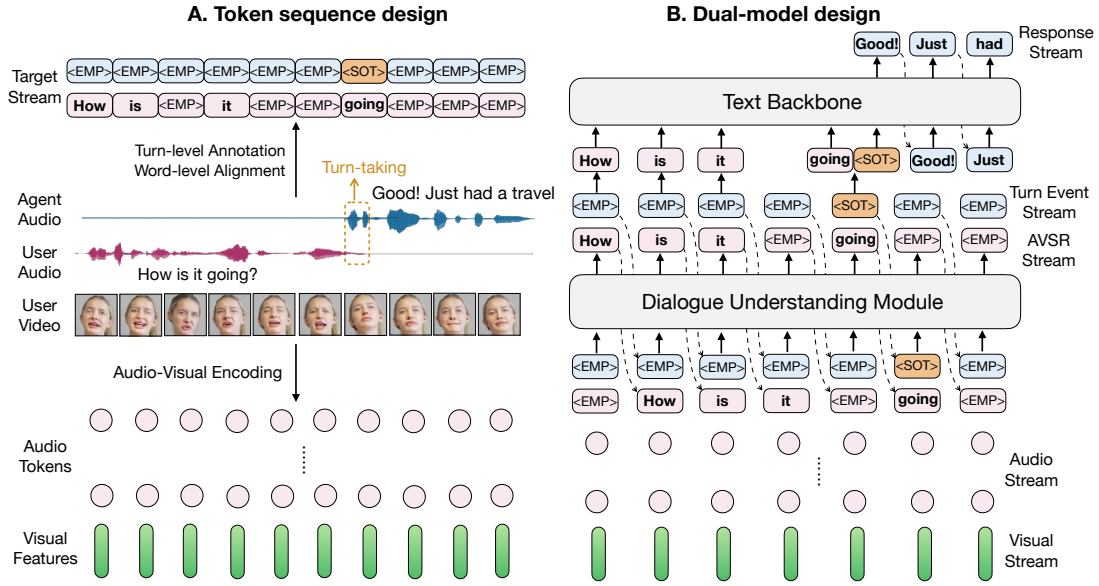


Figure 2: Token sequence and dual-model design. **A.** We use a DAC tokenizer to encode audio into 16 audio token streams and use AV-HuBERT to convert video to continuous visual features. We use turn-level annotation and word-level alignment to generate the target output text stream. **B.** shows our dual-model pipeline for AV-dialog. The AV dialogue understanding module recognizes user speech and detects potential turn-taking events, while the text backbone generates high-quality responses once turn-taking is triggered.

after the user finishes; (2) Overlapping turn, agent begins before user finishes, i.e. partial overlap; (3) Backchannel, short interjections (e.g., “hmm”, “yeah”). The model predicts special token <SOT> for both Normal and Overlapping turns, and <SOB> for Backchannels. For timesteps without turn taking events, we output <EMP>. An example turn-taking event stream is in Fig 2A.

3.2 Response Generation

Our audio-visual dialogue model, in its dual-model mode, streams user speech recognition and turn-taking predictions directly into a text-based LLM backbone (Fig. 2B). This design combines the responsiveness of instruction-tuned LLMs with the flexibility to swap in different text LLMs or APIs.

The text backbone operates in two states: *LISTENING* and *SPEAKING*. In *LISTENING*, non-silence tokens from the AV understanding module are streamed as the user’s input. When a turn-taking token appears, the model switches to *SPEAKING*, generating responses autoregressively. If new user speech tokens arrive mid-response, the model yields the floor and re-enters *LISTENING*.

To make response more natural and human-like, we explore two methods:

- **In-Context Learning (ICL):** We apply in-context learning (Brown et al., 2020) and add few-shot dialogue examples from SEAMLESS INTERACTION (InterAct) (Agrawal et al., 2025) training sets to the text backbone’s system prompt (see

§C.1.1).

- **Instruction Tuning (IT):** We finetune a chat-oriented LLM on real human dialogues using instruction tuning (Ouyang et al., 2022) to improve naturalness and responsiveness. The finetune hyperparameters can be found in (see §C.1.2).

The generated text is then converted to speech via the streaming TTS module Mimi (Défossez et al., 2024).

3.3 Training Strategy

A key challenge in training our AV-dialog understanding module is the scarcity of large-scale aligned audio-visual conversational data. To address this, we leverage diverse data sources: text datasets, monadic audio/audio-visual data, and real dyadic audio-only and audio-visual conversations. We build on a pre-trained text LLM, LLAMA-3-8B, and employ a two-stage, multi-task training approach to progressively develop the audio-visual understanding needed for a robust dialogue model.

3.3.1 Stage 1: Audio-Visual Understanding

The first stage focuses on aligning text, audio, and visual modalities through four multi-task objectives (see §B.1 for training hyper-parameters):

- *Text continuation:* Utilize large-scale text-only datasets for text continuation pre-training objective, to preserve robust language understanding and avoid catastrophic forgetting of textual data.

• *Speech comprehension*: Train on monaural speech datasets, LibriLight (Kahn et al., 2020), MLS (Pratap et al., 2020), and VP400k (Wang et al., 2021), on the ASR task to provide the model with acoustic comprehension of human speech.

• *Audio captioning*: Use the large audio dataset, Audioset (Gemmeke et al., 2017), for audio captioning task to achieve general audio comprehension.

• *Audio-visual alignment*: Train AVSR task using VoxCeleb2 (Nagrani et al., 2017), which is a large audio-visual monadic dataset. This enables the model to learn visual features linked to speech, fostering multimodal learning across text, audio, and visual modalities. To further improve AV understanding in noisy conditions, we apply the synthetic mixing augmentation from §3.3.3 on this dataset.

3.3.2 Stage 2: Learning about Conversations

We train the model on audio-only and audio-visual conversational data to learn natural dialogue dynamics. We use Fisher (Cieri et al., 2004) for audio-only and InterAct (Agrawal et al., 2025) for audio-visual conversations, optimizing two tasks: (1) streaming AVSR and (2) turn-taking event prediction. Synthetic mixing augmentation (§3.3.3) is also applied for robustness in noisy settings. Training hyperparameters are detailed in §B.2.

To prepare target sequences, we align words and turns by converting conversations into synchronized token streams. We deploy Whisper-Large (Radford et al., 2022) to acquire word-level timestamps. The first token of each word is placed at the $\lceil t_{start}/25 \rceil + d$ token in the AVSR stream, where t_{start} is the start timestamp from Whisper. The special turn event token is placed at $\lceil t_{turn}/25 \rceil$, where t_{turn} is the timestamp of the annotated turn event. Note that, the $d = 1s$ is also added to the AVSR stream but not to the Turn event stream to avoid introducing additional delays to the response.

3.3.3 Synthetic Mixing Augmentation

We apply synthetic mixing to simulate noisy, multi-speaker environments. For each training sample, with 20% probability, we use clean audio as input, with 40% probability, we mix the clean audio with background noise randomly sampled from from MUSAN (Snyder et al., 2015), and with 40% probability, we mix with 1–4 interference speakers from the same dataset. The input SNR is uniformly sampled between -8 dB and 8 dB. This augmentation enables the model to understand audio-visual cues in complex, real-world conditions.

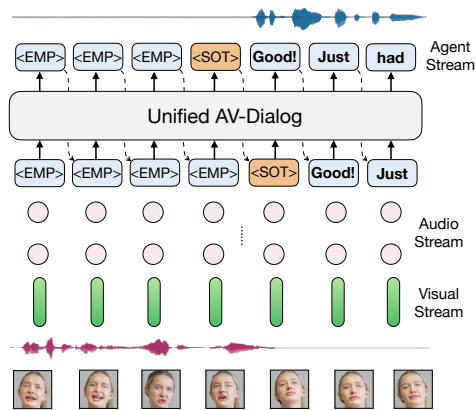


Figure 3: Unified AV-Dialog model. It takes the audio-visual input and predicts the turn-taking events. When the special turn-taking token is generated, the AV-Dialog model generates the response on the same output stream.

3.4 Unified AV-Dialog Model

We also explore a unified model variant where the AV-Dialog module directly generates full-duplex responses, eliminating the need for a text backbone (see §A for algorithmic latency analysis). This requires two key modifications:

• *Model and token design*: We remove the AVSR stream and instead add time-aligned agent response tokens interleaved with the turn-taking event stream (Fig. 3). Empirically, we found removal of AVSR stream improves unified model performance, enabling it to predict turn-taking events and then generate responses directly.

• *Training setup*: In Stage 2, we train on the Fisher and InterAct datasets to generate aligned text responses alongside turn-taking predictions (Fig. 3). Training hyperparameters are provided in §B.3.

4 Experiments

For each AV dialog sample, one side is randomly chosen as the user, whose audio and visual tokens are streamed into the model while output tokens are generated simultaneously. We focus on three aspects: (1) how well the model understands audio-visual input, (2) how accurately it predicts turn-taking, and (3) the quality of its responses.

To test robustness in noisy conditions, we evaluate under three conditions:

- *Clean*: Clean raw audio as input.
- *BG*: Clean raw audio mixed with background noise (music, chatter) from MUSAN (Snyder et al., 2015), input SNR range is -8 dB to 12 dB.
- *Interf*: Clean raw audio mixed with 1-4 interfering speakers from the same dataset as the target

WER(%) on Voxceleb2 ↓	Clean	BG	Interf
Auto-AVSR	26.8	48.2	71.8
Ours	17.4	35.6	38.8
WER(%) on LRS2 ↓	Clean	BG	Interf
Auto-AVSR	15.8	34.0	60.0
Ours	9.53	24.0	28.4

Table 1: Benchmarking streaming AVSR on the test set of Voxceleb2 and LRS2. We compare WER (%) between AUTO-AVSR and our AV-Dialog model.

WER(%)↓	Clean	BG	Interf
Auto-AVSR	32.2	60.4	93.0
Ours (A)	28.6	68.0	92.2
Ours (V)	67.8	67.8	67.8
Ours (A+V)	16.3	37.4	30.8

Table 2: Streaming AVSR on the InterAct test set. Ours (A): trained and test on audio-only input. Ours (V): trained and tested on visual-only input. Ours (A+V): trained and tested on audio-visual input.

speaker, input SNR range is -8dB to 12dB.

4.1 Audio-Visual Understanding Evaluation

We evaluate streaming AVSR using word error rate (WER). We compare our model with the state-of-the-art streaming AVSR model, Auto-AVSR (Ma et al., 2023). More recent AVSR works (Rouditchenko et al., 2024; Cappellazzo et al., 2025a) focus on offline settings, where models process the full recording before inference, which is an easier task. So, we benchmark against Auto-AVSR on Voxceleb2 and LRS2 dataset.

For a fair comparison, we first benchmark on the Voxceleb2 test set, as both our model and Auto-AVSR are trained on its training set. Since Voxceleb2 lacks text labels, we use Whisper-Large to transcribe clean speech as ground truth for WER. Table 1 compares audio-only, video-only, and audio-visual models. AV-Dialog achieves consistently lower WER than Auto-AVSR, with audio-visual input yielding the best performance. In Table 1, we also compare our model with Auto-AVSR on the LRS2 test set with manually labeled transcripts. Auto-AVSR is trained on LRS2, while our model is not, demonstrating out-of-domain generalization of our model’s audio-visual understanding.

We also evaluate our models for the streaming AVSR task on the test-set of the InterAct dataset. We use the transcription from InterAct as our ground-truth text to compute WER. As shown in Table. 2, our AV-dialogue model achieves much lower WER than audio-only or visual-only input,

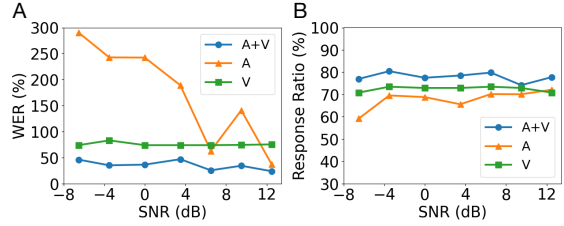


Figure 4: Model performance across different SNRs. Plot A shows the WER of streaming AVSR task on different SNRs of noisy audio input. Plot B shows the response ratio of the turn-taking prediction on different SNRs of noisy audio input.

demonstrating that combining modalities greatly improves AV understanding, especially under noise and interference. To further assess robustness, we evaluate across SNR ranges (Fig. 4A), averaging WER over multiple noise samples for BG and Interf scenarios. The results show our AV model remains robust across varying SNR levels.

4.2 Turn-Taking Prediction Evaluation

We evaluate our model using turn-taking events from (Nguyen et al., 2023) and measure floor-transfer offset (FTO), which is the duration between turn transitions, where negative FTO indicates overlap and positive FTO indicates a gap.

We extract the agent’s turn start timestamp using the SOT token and compute FTO as the gap between the user’s turn end (obtained from ground-truth turn annotations) and the agent’s turn start.

Metrics. We report three metrics:

- *Response Ratio*: percentage of FTOs within $-2s$ to $3s$, the typical range in InterAct’s human conversations (around 90% of FTOs in the InterAct test set are in this range).
- *FTO Error*: mean absolute error (MAE) between generated and ground-truth FTOs.
- *Median FTO*: the median value of FTOs.

Baselines. We compare with Moshi (Défossez et al., 2024), a state-of-art spoken dialogue model.

- **Moshi**: We deploy the Moshi checkpoints (Défossez et al., 2024) to generate responses.
- **SE+Moshi**: We first apply the speech enhancement (SE) model Demucs (Defossez et al., 2020), then feed output to Mosh to generate responses.
- **Separation+Moshi**: We first apply X-TF-GridNet (Hao et al., 2024), a state-of-the-art target speaker separation model, to the noisy mixture. We then feed output to Moshi to generate responses.

Results. As shown in Table 3, the audio-visual

Model	Clean			BG			Interf		
Metrics	Response Ratio(↑)	FTO Err(↓)	Median FTO	Response Ratio(↑)	FTO Err(↓)	Median FTO	Response Ratio(↑)	FTO Err(↓)	Median FTO
Moshi	54.0%	3.48	-0.72	53.8%	3.20	-0.12	52.5%	3.27	-0.52
SE + Moshi	55.9%	3.66	-0.6	56.0%	3.24	-0.7	54.4%	3.36	-0.84
Separation + Moshi	47.8%	3.60	-0.92	51.7%	3.47	-0.42	52.7%	3.27	-0.92
Ours (A)	73.2%	1.87	1.04	70.2%	2.66	1.02	65.8%	2.81	1.2
Ours (V)	72.5%	2.37	0.98	72.5%	2.37	0.98	72.5%	2.37	0.98
Ours (A+V)	74.5%	1.86	1.16	78.3%	1.96	1.12	78.8%	1.81	1.18
Ours (Unified)	68.1%	1.68	1.92	75.6%	1.49	1.76	75.9%	1.67	1.76
GT	-	0	1.5	-	0	1.5	-	0	1.5

Table 3: Turn-taking evaluation under different noise and interference conditions.

Noise condition	Clean		BG		Interf	
Model	PPL(↓)	Pickup Ratio(↑)	PPL(↓)	Pickup Ratio(↑)	PPL(↓)	Pickup Ratio(↑)
Moshi	44.1	23.8%	52.4	19.4%	46.4	18.4%
SE + Moshi	50.8	26.4%	50.4	24.5%	56.1	19.1%
Ours (ICL)	25.8	66.6%	24.6	68.1%	23.1	67.8%
Ours (IT)	23.2	32.5%	24.0	30.3%	23.8	36.7%
Ours (Unified)	31.0	29.6%	29.8	35.5%	32.5	31.3%
GT	51.7	-	51.7	-	51.7	-

Table 4: Semantic evaluation of dialogue responses under different noise conditions.

dialogue model achieves the highest response ratio across all noisy scenarios: 74.5% (*Clean*), 78.3% (*BG*), and 78.8% (*Interf*), substantially outperforming the baseline Moshi model, which reaches only 50%. Adding visual input to the audio-only model improves turn-taking accuracy by 1.3% (*Clean*), 8.1% (*BG*), and 13% (*Interf*). The unified audio-visual model shows a slight drop in response ratio but achieves the lowest FTO errors. Note that the GT median FTO is around 1.5s for the InterAct conversation dataset which consists of casual conversations between strangers. The detailed FTO distribution visualization can be found in §E.

To assess robustness under noise, in Fig. 4B, we also evaluate turn-taking prediction across different SNR ranges for both *BG* and *Interf* scenarios.

4.3 Semantic Evaluation

We evaluate the semantic quality of dialogue responses, comparing the Moshi baselines with three audio-visual dialogue variants:

- **Ours (ICL):** Dual-model pipeline using in-context learning with example conversations from InterAct in the prompt of the LLAMA3-8B text backbone model (see §C.1).
- **Ours (IT):** Dual-model pipeline fine-tuned via instruction tuning on Fisher, and InterAct datasets for the text backbone model (see §C.2).
- **Ours (Unified):** Unified model trained to generate text responses from audio-visual input.

We run our models end-to-end and compute the perplexity (PPL) of agent-generated turns. To further assess text quality, we use the Prometheus (Kim et al., 2023) LLM as an evaluator framework, performing relative/pairwise comparisons rather than absolute scoring, which better aligns with human judgment (Kiritchenko and Mohammad, 2017; Liusie et al., 2023). For each evaluation, the LLM compares the ground-truth InterAct response with the model-generated text. We compute the *Pickup Ratio* as the fraction of responses in which the LLM prefers the model-generated text over the ground truth (see details in §D).

As shown in Table 4, the baseline Moshi and SE+Moshi models achieve the lowest Pickup Ratio according to the LLM evaluator. The audio-visual dialogue model using in-context learning (ICL) achieves the highest Pickup Ratio among all methods. In contrast, the same dual-model pipeline fine-tuned via instruction tuning (IT) and the unified audio-visual model show a reduced Pickup Ratio of 30–40%. This drop is likely because InterAct dialogues often contain casual, unpredictable conversation; fine-tuning the generation task on such data can degrade response quality. For perplexity (PPL), both the ICL and IT audio-visual models achieve the lowest values. Finally, the response quality for the unified model is worse than the cascaded model settings for our AV-dialogue models. This is line with recent observations in the related domain of speech-to-speech dialog models (Hu

Model	N-MOS(↑)	H-MOS(↑)
SE + Moshi	2.39	2.10
Ours (Dual+ICL)	4.14	4.09
Ours (Unified)	3.54	3.02
GT	3.92	3.62

Table 5: Human evaluation result. N-MOS: Mean Opinion Score on Naturalness of response. H-MOS: Mean Opinion Score on Helpfulness of response.

et al., 2025b), where cascade model responses outperform unified models.

4.4 Human Evaluation

We conducted a human evaluation with 18 participants to assess the end-to-end performance of our audio-visual dialogue model. Model text outputs were converted to speech using the Moshi streaming TTS (Défossez et al., 2024), ensuring a fair comparison since both systems used the same TTS. Participants were given dialogue transcripts and audio, including both user turns and model responses.

We randomly selected 15 samples from the InterAct test set across Clean, BG, and Interf conditions (details and SNRs in §G). For each sample, participants evaluated four conditions: (1) SE+Moshi, (2) dual-model + ICL, (3) unified model, and (4) ground truth. Each participant rated 8 dialogue sets, with randomized method order to avoid bias. Ratings followed the Mean Opinion Score (MOS) protocol (ITU-T P.808 (ITU-T, 2018)) on a 5-point Likert scale, evaluating Naturalness (N-MOS) and Helpfulness (H-MOS) (see §F).

Table 5 shows that both our dual and unified models outperform the SE+Moshi baseline in naturalness and helpfulness. The dual model with in-context learning achieves the best results. The unified model drops by 0.5 in naturalness and 1.02 in helpfulness compared to the dual model. This is likely due to (1) the limited size of real conversational data and (2) the fact that the conversations in the dataset are mostly casual chit-chat, often containing low-quality responses, random topic shifts, and limited logical reasoning.

These results highlight that the dual model benefits from using real-world data primarily for turn-taking modeling while leveraging the pretrained text backbone for stronger generation quality. Notably, the MOS trends align with the LLM evaluator pick ratios in our semantic evaluation.

4.5 Ablation studies

We first compare acoustic tokens with semantic tokens. Using the same training setup, we trained the

AVSR WER(%) ↓	Clean	BG	Interf
DinoSR (A)	24.9	89.0	239.2
DinoSR (A+V)	26.9	83.0	67.0
Acoustic (A)	28.6	60.0	63.4
Acoustic (A+V)	16.3	37.4	30.8
Response Ratio(%) ↑	Clean	BG	Interf
DinoSR (A)	67.3	63.3	63.2
DinoSR (A+V)	69.5	49.1	47.8
Acoustic (A)	73.2	70.2	65.8
Acoustic (A+V)	74.5	76.9	78.8

Table 6: Acoustic & semantic token comparison.

AVSR WER(%) ↓	Clean	BG	Interf
Ours(A+V)	16.3	37.4	30.8
Ours(No Stage 1)	58.9	95.1	86.7
Ours(No Audio Diag)	22.6	37.8	31.8
Ours (No augmentation)	17.5	121.3	160.2

Table 7: Ablation Study on the training recipe.

AV model with DinoSR (Liu et al., 2023) semantic tokens instead of DAC tokens, as DinoSR is a newer, improved semantic representation compared to HuBERT. Table 6 shows that acoustic tokens outperform semantic tokens in both streaming AVSR and turn-taking prediction tasks.

Next, we compare different training strategies (Table 7). *No Stage 1* indicates skipping Stage 1 while training on the LLaMA3-8B model, while *No Audio Dialogue* excludes the audio-only dialogue dataset during Stage 2 fine-tuning. Including the audio-only dialogue dataset improves performance. *No augmentation* trains the second stage without Synthetic Mixing Augmentation. Results show a significant drop without Stage 1 or without Synthetic Mixing Augmentation.

We also evaluate the impact of explicit turn-taking supervision strategy on the unified model in Appendix. §J. Finally, we conduct ablation studies on the effect of different audio channels of tokenizer and visual input distortion in Appendix. §J.

5 Conclusion

We introduced AV-Dialog, the first streaming audio-visual dialogue system that integrates audio, vision, turn-taking, and response generation. Using acoustic tokenization, multi-stage training, and explicit turn-event supervision, it achieves robust performance in noisy, multi-speaker environments, outperforming audio-only baselines in transcription, turn-taking, and response quality. Human evaluations confirm that AV-Dialog enables more natural, helpful, and speaker-aware conversations, underscoring the value of combining listening and looking for real-world multimodal dialogue.

6 Limitations and Risks

Limitations. While AV-Dialog advances full-duplex dialogue in noisy environments, its performance can be further improved. It currently does not explicitly model non-verbal auditory cues (e.g., laughter, sighs) or visual cues (e.g., facial expressions, gestures) beyond lip movements. Enhancing the understanding and generation of these multimodal signals could make interactions more human-like. Finally, factors like poor lighting, occlusions (e.g., hands covering the mouth) or extreme head poses, can impair lip movement extraction (Shi et al., 2022), affecting speaker tracking and speech understanding. Developing lip encoders that are robust to such conditions is a promising and complementary direction for future work. This performance gap between cascaded and unified approaches is an important observation, and we would like to note that this has been highlighted as a limitation of unified approaches by prior work as well (Hu et al., 2025a; Veluri et al., 2024; Nguyen et al., 2023). Multiple prior works report a similar performance gap in response quality compared to cascaded baselines, with the primary driving factor being the lack of large-scale channel separated spoken dialog data.

Ethical considerations. Like any advanced AI enabling human-like interaction, AV-Dialog presents key ethical challenges. It may produce misleading dialogue, particularly under noisy or ambiguous conditions, requiring rigorous evaluation and ongoing monitoring. While audio-visual capture (e.g., lip movements, voices) is common in voice conferencing platforms like Zoom, it still demands strict attention to privacy. To prevent misuse such as exploitation in online scams, methods like speech watermarking could help safeguard against abuse.

References

Vasu Agrawal, Akinniyi Akinyemi, Kathryn Alvero, Morteza Behrooz, Julia Buffalini, Fabio Maria Carlucci, Joy Chen, Junming Chen, Zhang Chen, Shiyang Cheng, Praveen Chowdary, Joe Chuang, Antony D’Avirro, Jon Daly, Ning Dong, Mark Dupenthaler, Cynthia Gao, Jeffrey M. Girard, Martin Gleize, and 65 others. 2025. *Seamless interaction: Dyadic audiovisual motion modeling and large-scale dataset*. *CoRR*, abs/2506.22554.

Virginia Best, Alex D. Boyd, and Kamal Sen. 2023. *An effect of gaze direction in cocktail party listening*. *Trends in Hearing*, 27:23312165231152356. PMID: 36691678.

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

Umberto Cappellazzo, Minsu Kim, Honglie Chen, Pingchuan Ma, Stavros Petridis, Daniele Falavigna, Alessio Brutti, and Maja Pantic. 2025a. *Large language models are strong audio-visual speech recognition learners*. In *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5.

Umberto Cappellazzo, Minsu Kim, Honglie Chen, Pingchuan Ma, Stavros Petridis, Daniele Falavigna, Alessio Brutti, and Maja Pantic. 2025b. *Large language models are strong audio-visual speech recognition learners*. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.

Christopher Cieri, David Miller, and Kevin Walker. 2004. *The fisher corpus: a resource for the next generations of speech-to-text*. In *International Conference on Language Resources and Evaluation*.

Alexandre Defossez, Gabriel Synnaeve, and Yossi Adi. 2020. *Real time speech enhancement in the waveform domain*. *Preprint*, arXiv:2006.12847.

dlib. Dlib c++ library. <https://dlib.net/>.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. *The llama 3 herd of models*. *arXiv e-prints*, pages arXiv–2407.

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

Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing Moore, Manoj Plakal, and Marvin Ritter. 2017. *Audio set: An ontology and human-labeled dataset for audio events*. In *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 776–780. IEEE.

Fengyuan Hao, Xiaodong Li, and Chengshi Zheng. 2024. *X-tf-gridnet: A time–frequency domain target speaker extraction network with adaptive speaker embedding fusion*. *Information Fusion*, 112:102550.

Joanna Hong, Minsu Kim, Jeongsoo Choi, and Yong Man Ro. 2023. *Watch or Listen: Robust Audio-Visual Speech Recognition with Visual Corruption Modeling and Reliability Scoring*. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18783–18794, Los Alamitos, CA, USA. IEEE Computer Society.

722	Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai,	Zikai Liao, Yi Ouyang, Yi-Lun Lee, Chen-Ping Yu,	779
723	Kushal Lakhota, Ruslan Salakhutdinov, and Abdel-	Yi-Hsuan Tsai, and Zhaozheng Yin. 2025. Beyond	780
724	rahman Mohamed. 2021. Hubert: Self-supervised	words: Multimodal llm knows when to speak. <i>arXiv</i>	781
725	speech representation learning by masked prediction	<i>preprint arXiv:2505.14654</i> .	782
726	of hidden units . <i>Preprint</i> , arXiv:2106.07447.		
727	Ke Hu, Ehsan Hosseini-Asl, Chen Chen, Edresson	Alexander H Liu, Heng-Jui Chang, Michael Auli, Wei-	783
728	Casanova, Subhankar Ghosh, Piotr Żelasko, Zhehuai	Ning Hsu, and Jim Glass. 2023. Dinosr: Self-	784
729	Chen, Jason Li, Jagadeesh Balam, and Boris Gins-	distillation and online clustering for self-supervised	785
730	burg. 2025a. Efficient and direct duplex modeling	speech representation learning. <i>Advances in Neural</i>	786
731	for speech-to-speech language model. <i>arXiv preprint</i>	<i>Information Processing Systems</i> , 36:58346–58362.	787
732	<i>arXiv:2505.15670</i> .		
733	Ke Hu, Ehsan Hosseini-Asl, Chen Chen, Edresson	Adian Liusie, Potsawee Manakul, and Mark JF Gales.	788
734	Casanova, Subhankar Ghosh, Piotr Żelasko, Zhehuai	2023. Llm comparative assessment: Zero-shot nlg	789
735	Chen, Jason Li, Jagadeesh Balam, and Boris Gins-	evaluation through pairwise comparisons using large	790
736	burg. 2025b. Salm-duplex: Efficient and direct du-	language models. <i>arXiv preprint arXiv:2307.07889</i> .	791
737	plex modeling for speech-to-speech language model .		
738	<i>Preprint</i> , arXiv:2505.15670.	Pingchuan Ma, Alexandros Haliassos, Adriana	792
739	ITU-T. 2018. Recommendation P.808: Subjective eval-	Fernandez-Lopez, Honglie Chen, Stavros Petridis,	793
740	uation of speech quality with a crowdsourcing ap-	and Maja Pantic. 2023. Auto-avsr: Audio-visual	794
741	proach . ITU-T Recommendation P.808, International	speech recognition with automatic labels. In <i>ICASSP</i>	795
742	Telecommunication Union.	2023-2023 <i>IEEE International Conference on Acous-</i>	796
743	J. Kahn, M. Rivière, W. Zheng, E. Kharitonov, Q. Xu,	<i>tics, Speech and Signal Processing (ICASSP)</i> , pages	797
744	P. E. Mazaré, J. Karadayi, V. Liptchinsky, R. Col-	1–5. IEEE.	798
745	lobert, C. Fuegen, T. Likhomanenko, G. Synnaeve,	T Aleksandra Ma, Sile Yin, Li-Chia Yang, and Shuo	799
746	A. Joulin, A. Mohamed, and E. Dupoux. 2020. Libri-	Zhang. 2025. Real-time audio-visual speech en-	800
747	light: A benchmark for asr with limited or no su-	hancement using pre-trained visual representations.	801
748	perception. In <i>ICASSP 2020 - 2020 IEEE Interna-</i>	<i>arXiv preprint arXiv:2507.21448</i> .	802
749	<i>tional Conference on Acoustics, Speech and Signal</i>	Josh Mcdermott. 2009. The cocktail party problem .	803
750	<i>Processing (ICASSP)</i> , pages 7669–7673. https://github.com/facebookresearch/libri-light .	<i>Current biology : CB</i> , 19:R1024–7.	804
751		Arsha Nagrani, Joon Son Chung, and Andrew Zisser-	805
752	Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang,	man. 2017. Voxceleb: a large-scale speaker identifi-	806
753	Shayne Longpre, Hwaran Lee, Sangdoon Yun,	cation dataset . <i>CoRR</i> , abs/1706.08612.	807
754	Seongjin Shin, Sungdong Kim, James Thorne, and 1	Tu Anh Nguyen, Eugene Kharitonov, Jade Copet, Yossi	808
755	others. 2023. Prometheus: Inducing fine-grained	Adi, Wei-Ning Hsu, Ali Elkahky, Paden Tomasello,	809
756	evaluation capability in language models. <i>arXiv</i>	Robin Algayres, Benoit Sagot, Abdelrahman Mo-	810
757	<i>preprint arXiv:2310.08491</i> .	hamed, and Emmanuel Dupoux. 2022. Genera-	811
758	Svetlana Kiritchenko and Saif M Mohammad. 2017.	tive spoken dialogue language modeling . <i>Preprint</i> ,	812
759	Best-worst scaling more reliable than rating scales:	arXiv:2203.16502.	813
760	A case study on sentiment intensity annotation. <i>arXiv</i>	Tu Anh Nguyen, Eugene Kharitonov, Jade Copet, Yossi	814
761	<i>preprint arXiv:1712.01765</i> .	Adi, Wei-Ning Hsu, Ali Elkahky, Paden Tomasello,	815
762	Rithesh Kumar, Prem Seetharaman, Alejandro Luebs,	Robin Algayres, Benoit Sagot, Abdelrahman Mo-	816
763	Ishaan Kumar, and Kundan Kumar. 2023. High-	hamed, and 1 others. 2023. Generative spoken dia-	817
764	fidelity audio compression with improved rvqgan.	logue language modeling. <i>Transactions of the Asso-</i>	818
765	In <i>Proceedings of the 37th International Conference</i>	<i>ciation for Computational Linguistics</i> , 11:250–266.	819
766	<i>on Neural Information Processing Systems, NIPS '23</i> ,	Tu Anh Nguyen, Benjamin Muller, Bokai Yu, Marta R.	820
767	Red Hook, NY, USA. Curran Associates Inc.	Costa-jussa, Maha Elbayad, Sravya Popuri, Paul-	821
768	Kushal Lakhota, Evgeny Kharitonov, Wei-Ning Hsu,	Ambroise Duquenne, Robin Algayres, Ruslan Mav-	822
769	Yossi Adi, Adam Polyak, Benjamin Bolte, Tu-Anh	lyutov, Itai Gat, Gabriel Synnaeve, Juan Pino, Benoit	823
770	Nguyen, Jade Copet, Alexei Baevski, Adelrahman	Sagot, and Emmanuel Dupoux. 2024. Spirit-lm:	824
771	Mohamed, and Emmanuel Dupoux. 2021. Genera-	Interleaved spoken and written language model .	825
772	tive spoken language modeling from raw audio .	<i>Preprint</i> , arXiv:2402.05755.	826
773	<i>Preprint</i> , arXiv:2102.01192.	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	827
774	Sean Leishman, Peter Bell, and Sarenne Wallbridge.	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	828
775	2024. Pairwiseturngpt: a multi-stream turn predic-	Sandhini Agarwal, Katarina Slama, Alex Ray, and 1	829
776	tion model for spoken dialogue. In <i>Proceedings of</i>	others. 2022. Training language models to follow in-	830
777	<i>the 28th Workshop on the Semantics and Pragmatics</i>	structions with human feedback. <i>Advances in neural</i>	831
778	<i>of Dialogue</i> .	<i>information processing systems</i> , 35:27730–27744.	832

833	Se Park, Chae Kim, Hyeongseop Rha, Minsu Kim,	intrinsic cross-modal conversational abilities. In	890
834	Joanna Hong, Jeonghun Yeo, and Yong Ro. 2024.	<i>Findings of the Association for Computational Lin-</i>	891
835	Let’s go real talk: Spoken dialogue model for face-	<i>guistics: EMNLP 2023, Singapore, December 6-10,</i>	892
836	to-face conversation . In <i>Proceedings of the 62nd An-</i>	<i>2023</i> , pages 15757–15773. Association for Computa-	893
837	<i>annual Meeting of the Association for Computational</i>	<i>tional Linguistics.</i>	894
838	<i>Linguistics (Volume 1: Long Papers)</i> , pages 16334–		
839	16348, Bangkok, Thailand. Association for Computa-		
840	tional Linguistics.		
841	Vineel Pratap, Qiantong Xu, Anuroop Sriram, Gabriel	A Algorithm Latency Analysis	895
842	Synnaeve, and Ronan Collobert. 2020. MIs: A large-	Moshi is built on a 7B backbone model, while our	896
843	scale multilingual dataset for speech research. <i>ArXiv</i> ,	unified system uses a single 8B LLaMA model.	897
844	abs/2012.03411 .	Our dual-model architecture extends this by run-	898
845	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brock-	ning two such models in parallel, which naturally	899
846	man, Christine McLeavey, and Ilya Sutskever. 2022.	increases peak memory consumption.	900
847	Robust speech recognition via large-scale weak su-	Because system latency depends on hardware	901
848	pervision . <i>Preprint</i> , arXiv:2212.04356.	and software optimizations, which is not the fo-	902
849	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brock-	cus of our paper, we focus instead on algorithmic	903
850	man, Christine McLeavey, and Ilya Sutskever. 2023.	latency which is a platform-invariant metric. In	904
851	Robust speech recognition via large-scale weak super-	AV-Dialog, both the audio tokenizer (DAC) and	905
852	vision . In <i>Proceedings of the 40th International Con-</i>	the visual encoder operate causally at 25 Hz. How-	906
853	<i>ference on Machine Learning, ICML’23</i> . JMLR.org.	ever, the AV-HuBERT visual encoder introduces	907
854	Andrew Rouditchenko, Yuan Gong, Samuel Thomas,	a 2-frame lookahead (Ma et al., 2025), result-	908
855	Leonid Karlinsky, Hilde Kuehne, Rogerio Feris, and	ing in an overall algorithmic latency of approxi-	909
856	James Glass. 2024. Whisper-flamingo: Integrating	120 ms. In our dual-model setup, output tokens	910
857	visual features into whisper for audio-visual speech	from the understanding module are streamed di-	911
858	recognition and translation . In <i>Interspeech 2024</i> ,	rectly to the text backbone with KV-cache in par-	912
859	pages 2420–2424.	allel, enabling immediate response generation dur-	913
860	Bowen Shi, Wei-Ning Hsu, Kushal Lakhotia, and Ab-	turn-taking without additional delay.	914
861	delrahman Mohamed. 2022. Learning audio-visual	By contrast, Moshi (Défossez et al., 2024) pro-	915
862	speech representation by masked multimodal cluster	cesses 80 ms audio chunks, achieving an algori-	916
863	prediction . In <i>International Conference on Learning</i>	thmic latency of about 80 ms. Note that AV-Dialog	917
864	<i>Representations</i> .	employs Moshi’s TTS model. Although our sys-	918
865	David Snyder, Guoguo Chen, and Daniel Povey. 2015.	tem’s latency is somewhat higher, it is mainly lim-	919
866	Musan: A music, speech, and noise corpus . <i>arXiv</i>	ited by the visual encoder’s lookahead: an aspect	920
867	preprint arXiv:1510.08484 .	that could be further reduced by pretraining a vi-	921
868	Bandhav Veluri, Benjamin N Peloquin, Bokai Yu,	sual encoder with a smaller or zero lookahead win-	922
869	Hongyu Gong, and Shyamnath Gollakota. 2024. Be-	ow. Despite this limitation compared to Moshi-like	923
870	yond turn-based interfaces: Synchronous LLMs as	natively full-duplex models, we believe our approach	924
871	full-duplex dialogue agents . In <i>Proceedings of the</i>	of predicting agent’s start-of-the-turn bridges the	925
872	<i>2024 Conference on Empirical Methods in Natural</i>	naturalness gap while also leveraging superior help-	926
873	<i>Language Processing</i> , pages 21390–21402, Miami,	fulness and knowledge of standalone text back-	927
874	Florida, USA. Association for Computational Lin-	bones.	928
875	guistics.	B Training Hyper-parameters	929
876	Changhan Wang, Morgane Riviere, Ann Lee, Anne Wu,	B.1 Stage-1 Training	930
877	Chaitanya Talnikar, Daniel Haziza, Mary Williamson,	In Stage 1, we trained the original LLAMA3-8B	931
878	Juan Pino, and Emmanuel Dupoux. 2021. VoxPop-	with sequence length 4096. We use a learning	932
879	uli: A large-scale multilingual speech corpus for rep-	rate of $3e^{-5}$ on the transformer block and a learn-	933
880	resentation learning, semi-supervised learning and	ing rate of $1.5e^{-4}$ on embedding layers and au-	934
881	interpretation . In <i>Proceedings of the 59th Annual</i>	dio/visual adapters. The model is trained with 500	935
882	<i>Meeting of the Association for Computational Lin-</i>	step warmup and trained for 50k iterations on 128	936
883	<i>guistics and the 11th International Joint Conference</i>	A100 GPUs with a per-gpu batch size of 1.	937
884	<i>on Natural Language Processing (Volume 1: Long</i>		
885	<i>Papers)</i> , pages 993–1003, Online. Association for		
886	Computational Linguistics.		
887	Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan,		
888	Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. 2023.		
889	Speechgpt: Empowering large language models with		

The proportion of each task and dataset in the stage1 training is as follows:

- *Text continuation* (48.0%): Arxiv (16.0%), B3g (20.0%) and Wikipedia_en (12.0%).
- *Speech comprehension* (32.0%): LibriLigh (13.76%), MLS (13.12%) and VP400k (5.12%).
- *Audio captioning* (4.0%): AudioSet (4.0%).
- *Audio-visual alignment (AVSR)* (16.0%): Vox-celeb2 (16.0%).

The input and output token sequence design for Stage 1 training is shown in Fig. 5. The special tokens <ASR>, <Trans>, and <AC> serve as prefixes for different tasks. We also introduce a special <NULL> token: when a modality is missing in the input stream for a given task, it is filled with <NULL>, whose embedding vector is all zeros after the embedding layer. If <NULL> appears in the target stream, its loss is not computed. We apply cross-entropy loss on the output text stream.

B.2 Stage-2 Training for Dual Model

In Stage 2, we fine-tuned the Stage 1 model with a sequence length of 4096. We used a learning rate of $2e^{-5}$ for the transformer blocks, embedding layers, and audio/visual adapters. The model was trained with a 500-step warm-up over 10k iterations on 32 A100 GPUs, with a per-GPU batch size of 1. The proportion of each task and dataset in Stage 2 fine-tuning is as follows:

- *Audio-only conversation* (55.0%): Fisher (10.0%) and InterAct (45.0%).
- *Audio-Visual conversation* (45.0%): InterAct (45.0%).

The input and output token sequence design for Stage 2 training is shown in Fig. 6. We also apply the <NULL> token in the same way as Stage 1. We compute the cross-entropy loss on both the AVSR stream U and Turn event stream T and compute their average. In the AVSR stream, the loss weight for text tokens is set to 1.0, while the silence token <EMP> is set to 0.1. In the Turn-Event stream, the loss weight for Turn-taking token <SOT> is set to 2.5, the loss weight for the backchannel token <BOT> is set to 1.0, and the loss weight for the the silence token <EMP> is set to 0.1.

B.3 Stage-2 Training for Unified Model

In Stage 2 of our unified model, we fine-tuned the pretrained Stage 1 model with a sequence length of 4096. We used a learning rate of $2e^{-5}$ for the trans-

former blocks, embedding layers, and audio/visual adapters. The model was trained with a 500-step warm-up over 10k iterations on 32 A100 GPUs, with a per-GPU batch size of 1.

The proportion of each task and dataset in Stage 2 fine-tuning is as follows:

- *Audio-only conversation* (55.0%): Fisher (10.0%), InterAct (45.0%).
- *Audio-Visual conversation* (45.0%): InterAct (45.0%).

The input and output token sequence design for the unified model training at Stage 2 is shown in Fig. 7. We apply the <NULL> token in the same way as Stage 1. We compute the cross-entropy loss on the output stream T . The loss weight for text tokens is set to 1.0, the loss weight for <EMP> is set to 0.1, the loss weight for <SOT> is set to 2.5, and the loss weight for <BOT> is set to 1.0.

C Text Backbone Hyper-parameters

C.1 In-Context Learning

The prompting and few-shot samples are provided as:

"Carefully read the user prompt. You follow these instructions:

You are a helpful assistant that engages in natural, casual conversation. Respond like a human would - be conversational, use natural language, and don't be overly formal.

Here are some examples of the conversational style you should adopt:

Example1:

user: Hey, Siobhan, what's up? You seem troubled.

assistant: Yeah, I am. I'm just having a hard time. I needed someone to talk to.

user: Of course, man, I'm always here for you. What's going on?

assistant: It's just everything. Work is stressing me out. My relationship is falling apart, and I feel like I'm losing touch with my friends. I don't know what to do.

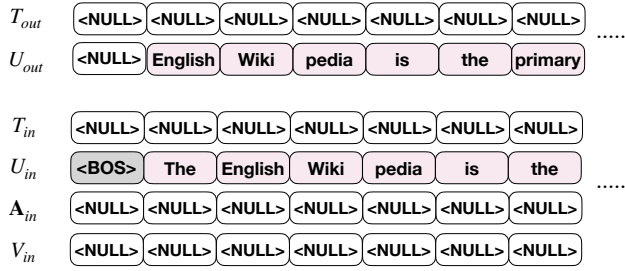
user: Well, let's start with work then. What's going on there?

assistant: It's just that... Everything is so demanding and I can't keep up. I'm constantly behind and it feels like I'm never gonna catch up.

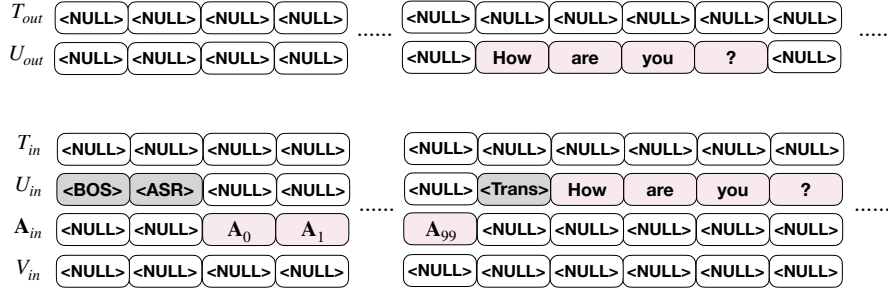
Example2:

user: Hey, thanks for taking my motorcycle off my hands. I really appreciate it.

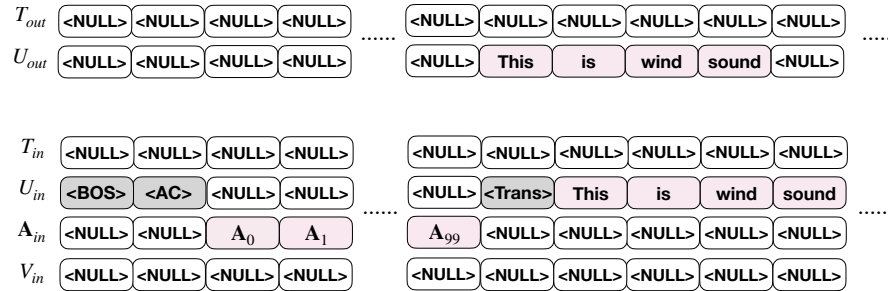
A. Text continuation task



B. Speech Comprehension



C. Audio Captioning



D. Audio-visual Alignment

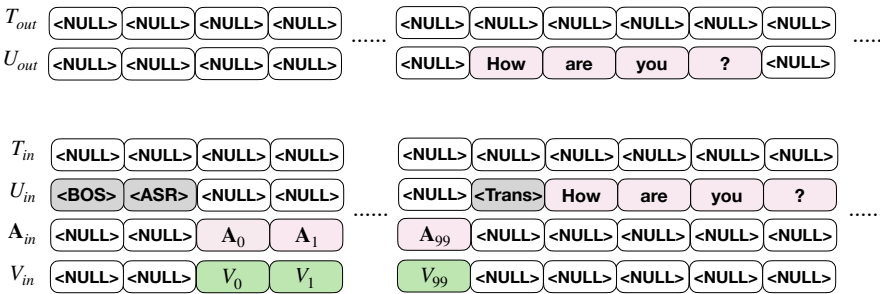


Figure 5: Training input and output tokens design for different tasks at Stage 1.

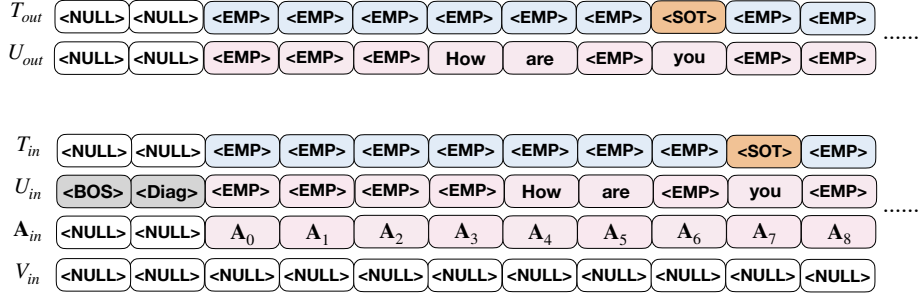
1035
1036
1037
1038
1039
1040
1041
1042
1043

assistant: No problem, man. It's a nice bike. I'm sure I'll enjoy riding it.
user: I'm just glad to be rid of it, to be honest. It was costing me a lot in repairs and maintenance.
assistant: Yeah, motorcycles can be pretty expensive to keep up. They're also a lot of fun to ride. So it's worth it in my opinion.
user: I'm not so sure. I think I'd rather just... take the bus from now on. It's cheaper and less

hassle.
assistant: That's a valid point. But sometimes it's nice to have the freedom that a motorcycle provides. You can go where you want, when you want.
user: I guess that's true, but it's just not worth the expense for me anymore.
Example3: user: Would you rather have the ability to speak to the past or send messages to the

1044
1045
1046
1047
1048
1049
1050
1051
1052

A. Audio-only Conversation



B. Audio-Visual Conversation

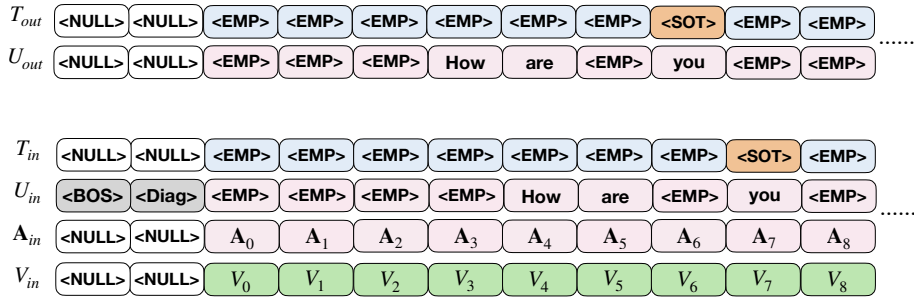
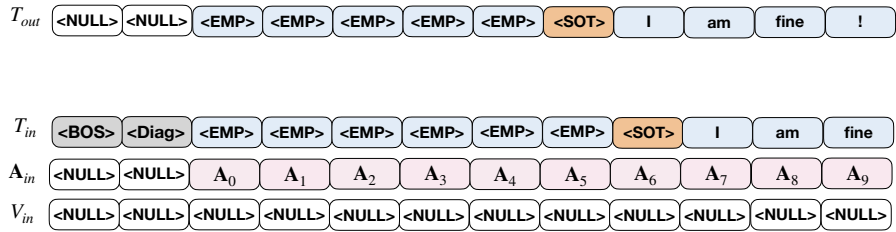


Figure 6: Training input and output tokens design for different tasks at Stage 2.

A. Audio-only Conversation for unified model



B. Audio-Visual Conversation for unified model

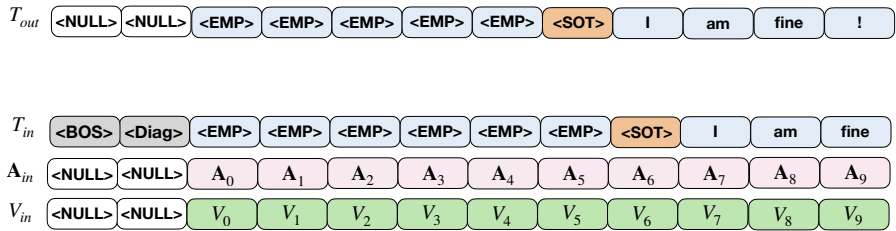


Figure 7: Training input and output tokens design for multi-tasks at Stage 2 for the unified model.

1053	future?		
1054	assistant: I would rather send messages to the		
1055	future.		
1056	user: Really? But we can kind of already do		
1057	that, don't you think, with certain apps and maybe		
1058	hiding a letter in your house or something?		
1059	assistant: Well, to me, the progress, I like the		
1060	progress, what's going on, and I can't just wait to		
1061	see what else up there.		
1062	user: What kind of messages would you want to		
1063	send?		
1064	assistant: And messages are positive, always		
1065	positive.		
1066	Now respond to the user in this same natural		
1067	conversational style. Remember your instructions."		
1068	C.2 Instruction Tuning		
1069	We finetune LLAMA3-8B-INSTRUCT with a se-		
1070	quence length of 8192 on human conversational		
1071	datasets, including 22.2% InterAct dataset, 3.7%		
1072	Fisher dataset. We pre-processing the dialogue		
1073	dataset by applying the original instruction-tuning		
1074	template from Llama3-8B-Instruct. We use a low		
1075	learning rate of $1e^{-5}$ and finetuned it for only 3000		
1076	steps on 32 A100 GPUs with a per-gpu batch size		
1077	of 1.		
1078	D Evaluation Prompting for Prometheus		
1079	We deploy PROMETHEUS-7B-V2.0 as our LLM		
1080	evaluator on the generated response. We randomly		
1081	shuffle the order of ground-truth and generated re-		
1082	sponse. The rubric description is as follow:		
1083	"Does Agent respond in a way that is generally		
1084	related to the user's input and current conversa-		
1085	tion? Minor topic drift, informal language, brevity,		
1086	or slight ambiguity should not be penalized. Ig-		
1087	nore formatting, punctuation, and minor inconsis-		
1088	tencies."		
1089	E FTO Distribution		
1090	Fig. 8 visualizes the different FTO distributions for		
1091	different model configurations:		
1092	• A. Moshi: state of the art speech-only dialogue		
1093	model		
1094	• B. SE+Moshi: cascade of a speech enhancement		
1095	model (Demcus) and Moshi		
1096	• C. Semantic(A): our dual-model approach with		
1097	semantic tokens (DinoSR) and audio-only input		
1098	• D. Semantic(A+V): our dual-model approach		
1099	with semantic tokens (DinoSR) and audio-visual		
	input		1100
	• E. Acoustic(A): our dual-model approach with		1101
	acoustic tokens (DAC) and audio-only input		1102
	• F. Acoustic(V): our dual-model approach with		1103
	acoustic tokens (DAC) and visual-only input		1104
	• G. Acoustic(A+V): our dual-model approach		1105
	with acoustic tokens (DAC) and audio-visual in-		1106
	put		1107
	• G. Unified(A+V): our unified-model approach		1108
	with acoustic tokens (DAC) and audio-visual input		1109
	F N-MOS and H-MOS		1110
	N-MOS scores for Naturalness are defined as fol-		1111
	lows:		1112
	• 1. Bad - Response is not normal English or does		1113
	not make sense.		1114
	• 2. Poor - Response is normal English but not		1115
	coherent to the user's input.		1116
	• 3. Fair - Response is somewhat plausible and		1117
	coherent		1118
	• 4. Good - Response is plausible and coherent		1119
	• 5. Excellent - Response is highly plausible and		1120
	coherent		1121
	M-MOS scores for Meaningfulness are defined		1122
	as follows:		1123
	• 1. Bad - essentially nothing in common with		1124
	human-like conversation		1125
	• 2. Poor - very little natural and human-like con-		1126
	versation		1127
	• 3. Fair - substantial differences from human-like		1128
	and natural conversation		1129
	• 4. Good - minor differences from human-like		1130
	and natural conversation		1131
	• 5. Excellent - basically indistinguishable from		1132
	human-like and natural conversation		1133
	G Samples distribution of real-human		1134
	evaluation		1135
	The properties of the samples used in human evalua-		1136
	tion are shown in Table. 8.		1137
	H User study participants		1138
	The human evaluation study was performed under		1139
	our institution's IRB. All participants provided con-		1140
	sent and were recruited from our institutions and		1141
	nearby areas. Participants ranged from 18 to 40		1142
	years old, with 35% identifying as female and the		1143

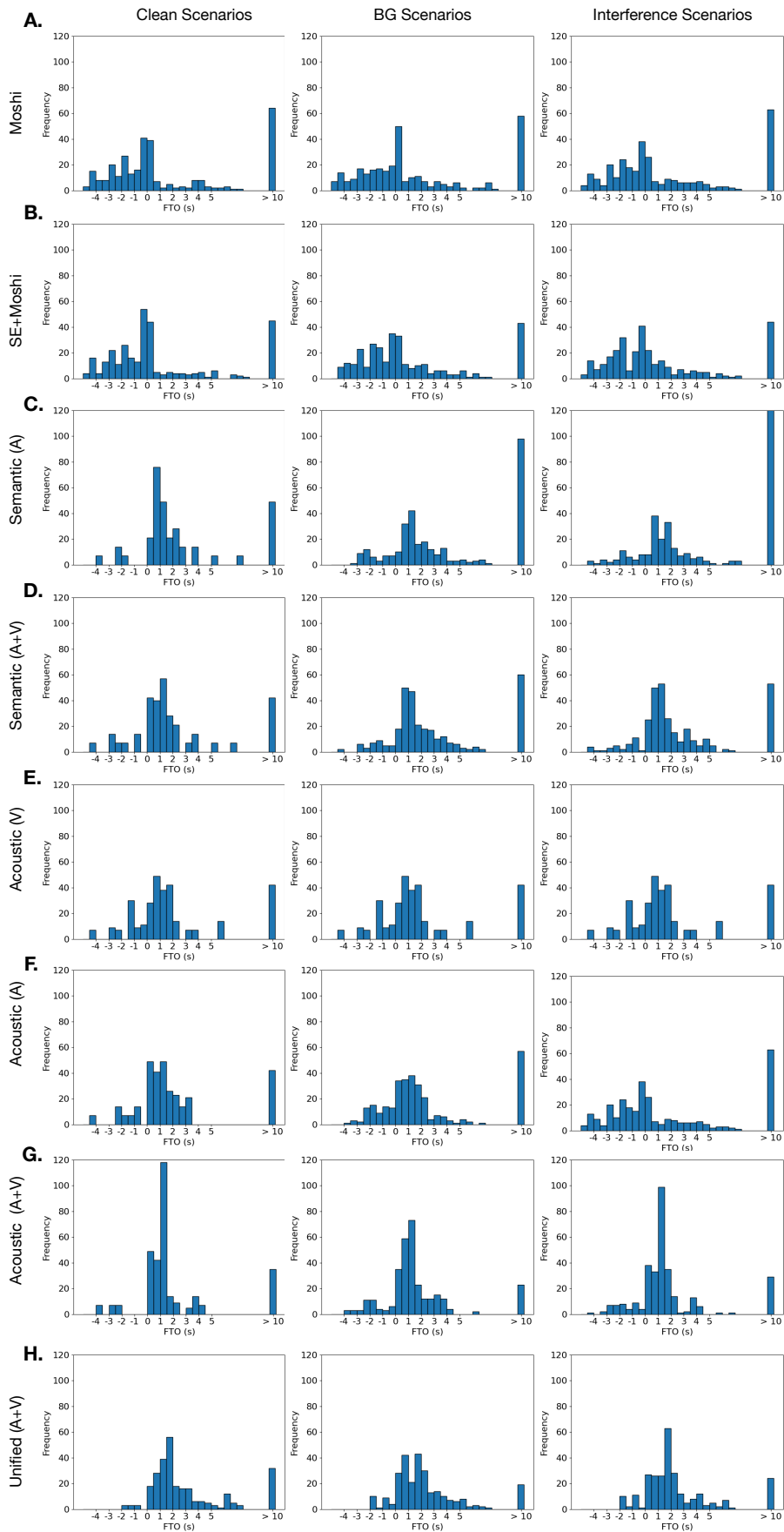


Figure 8: Distribution of FTO of different model configurations. On the x-axis of each figure, the label "> 10" also accounts for samples where the model does not respond at all.

Sample Index	Noise Scenario	SNR(dB)
1	BG	9.0
2	BG	0.0
3	clean	∞
4	Interf	3.0
5	clean	∞
6	Interf	-2.99
7	Interf	2.99
8	clean	∞
9	clean	∞
10	Interf	-3.0
11	BG	3.0
12	BG	6.0
13	Interf	-7.0
14	BG	3.0
15	BG	0.0

Table 8: Distribution of the samples used in human evaluation.

AVSR WER(%) ↓	Clean	BG	Interf
16 channels	16.3	37.4	30.8
8 channels	19.2	36.6	40.1
4 channels	23.8	65.0	48.8
Response Ratio(%) ↑	Clean	BG	Interf
16 channels	74.5	78.3	78.8
8 channels	79.1	78.4	77.1
4 channels	73.8	78.6	72.5

Table 9: Ablation Study on the audio channel number.

rest as male. They are recruited from both technical background and non-technical background.

I Example of AV-Dialog

We provide some samples of our AV-Dialog input and output in Figure. 9

J Ablation Study

J.1 Effect of audio channels

To evaluate the impact of audio compression quality, we experimented with varying the number of audio token channels. Since our Residual Vector Quantization (RVQ) audio tokenizer encodes information hierarchically (coarse-to-fine), we compared the full 16 channels against reduced inputs of 8 and 4 channels. As shown in Table. 9, reducing the input to 8 channels results in only a slight performance drop. Another interesting observation is that the AV understanding task is more vulnerable to the audio channel drop, because AV understanding task relies heavily on the fine-grained acoustic details captured in the deeper RVQ levels, which are lost when channels are aggressively pruned.

AVSR WER(%) ↓	Clean	BG	Interf
0% drop rate	16.3	37.4	30.8
5% drop rate	17.3	34.0	32.6
10% drop rate	18.4	35.8	33.5
25% drop rate	17.9	35.6	35.6

Table 10: Ablation Study on the visual frame drop.

Response ratio(%) ↑	Clean	BG	Interf
w explicit turn-taking	68.1	75.6	75.9
w/o explicit turn-change	48.0	35.1	38.0
LLM-evaluator pick-up ratio(%) ↑	Clean	BG	Interf
w explicit turn-change	29.6	35.5	31.3
w/o explicit turn-change	29.0	22.1	18.2

Table 11: Ablation study on turn-taking supervision in our Unified Model.

J.2 Effect of Visual Frame Drop

We conducted additional experiments on video modality distortion (e.g., occlusions, back-facing, or extreme angles). In our pipeline, such distortions cause the face detection module to fail, resulting in dropped visual frames. We simulated this by randomly dropping video frames at different rates. The results in Table. 10 show that our model can be robust to the video frame dropped due to the severe distortion.

J.3 Effect of Explicit Turn-taking Supervision

We also evaluate the impact of explicit turn-taking supervision on the unified model. In Stage 2, we trained a version without the <SOT> token, forcing it to generate response tokens directly. Table 11 shows that this supervision is crucial for both turn-taking prediction and response quality.

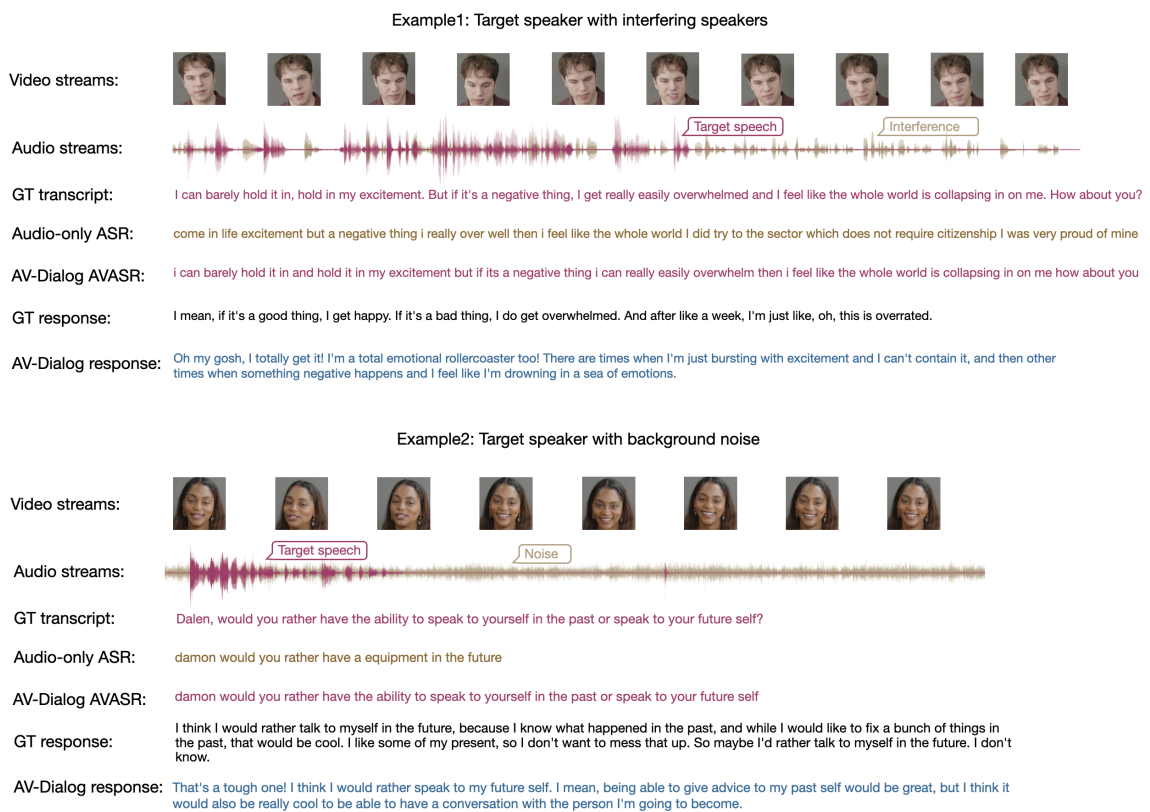


Figure 9: Data samples for our AV-dialog. Example 1 has interfering speakers while Example 2 has significant background noise. Video and audio streams are input to AV-Dialog. GT transcript is the ground-truth transcription of the target speaker. Audio-only ASR is the audio-only model output. AV-Dialog AVASR is the output of AVSR stream from AV-Dialog. GT response is the ground-truth response and AV-Dialog response is the output of our dialogue model.