

AudioRole: An Audio Dataset for Character Role-Playing in Large Language Models

Anonymous ACL submission

Abstract

While existing role-playing research predominantly focuses on text, **Audio Role-Playing (ARP)** presents unique challenges regarding the synchronized alignment of semantic content and vocal characteristics. To address this gap, we propose **AudioRole**, a meticulously curated dataset from 13 TV series spanning 1K+ hours with 1M+ character-grounded dialogues, providing synchronized audio-text pairs annotated with speaker identities and contextual metadata. In addition, to demonstrate the effectiveness of the dataset, we introduced **ARP-Eval**, a dual-aspect evaluation framework that assesses both *response quality* and *role fidelity*. Empirical validation showing GLM-4-Voice trained on AudioRole (called **ARP-Model**) achieves an average Acoustic Personalization score of 0.31, significantly outperforming the original GLM-4-voice and the more powerful model MiniCPM-O-2.6. The **ARP-Model** also achieves a Content Personalization score of 0.36, surpassing the untrained original model by about 38%. The blind human perceptual evaluation also confirms these findings.

AudioRole features dialogues from over 115 main characters, 6 trained **ARP-Models**, and evaluation protocols. Together, they provide an essential resource for advancing audio-grounded role-playing research.

1 Introduction

The evolution of role-playing capabilities in LLMs has revolutionized human-AI interaction paradigms. Contemporary systems demonstrate remarkable proficiency in textual persona simulation, serving as personalized assistants¹, emotional companions², and social interaction proxies (Park et al., 2023). However, this progress remains fundamentally constrained by unimodal (text-only) interaction frameworks, which neglect the critical

¹<https://chatgpt.com/>

²<https://replika.com/>

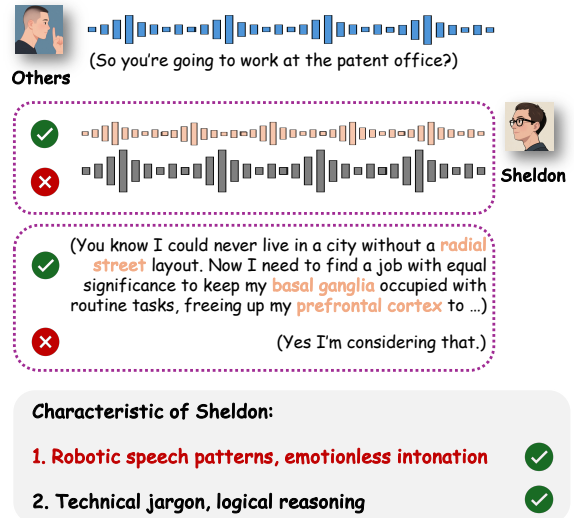


Figure 1: Audio Role-Playing case of Sheldon. The answer should satisfy not only the lexical similarity but also acoustic similarity.

role of vocal characteristics in authentic character portrayal — a limitation that is particularly evident when simulating personas from audiovisual media, such as films and TV series.

Existing work on audio evaluation considers mainly content accuracy (Chen et al., 2024b) (Hasid et al., 2024) or only the conversational style (Liu et al., 2025). Vocal expression reflects an integral component of communication that varies considerably between individuals in different contexts (Cohen et al., 2015), and the acoustic properties of speech reflect key variables for understanding human behavior (Decety and Lamm, 2006). Vocal expression is highly variable across individuals and contexts, and is influenced by several individual differences. Consider Dr. Sheldon Cooper from *The Big Bang Theory* (Figure 1): while text-based LLMs can mimic his technical jargon and logical reasoning, they fundamentally fail to capture his iconic robotic speech patterns — specifically, his faster-than-average speaking rate and emotionless

062 intonation. This disconnect highlights a critical
063 challenge in audio-grounded role-playing: authentic
064 characterization requires not just *what* the person
065 says, but precisely *how* they say it — the
066 synchronized alignment of semantic content and
067 acoustic delivery that defines recognizable person-
068 alities.

069 To promote the solution of this challenge, we
070 introduce **Audio Role-Playing (ARP)**, a novel task
071 requiring dual-alignment of generated responses:
072 (1) semantic consistency with character knowledge
073 and speaking style, and (2) acoustic fidelity to vo-
074 cal profiles. ARP task extends beyond speech syn-
075 thesis — successful ARP systems must dynami-
076 cally adapt both *what* is said (content) and *how*
077 it’s said (delivery) according to situational contexts
078 and character traits.

079 Central to advancing ARP research is the crea-
080 tion of high-quality training data. Current at-
081 tempts using general-purpose multi-modal mod-
082 els such as GPT-4o-Audio and MiniCPM-O-2.6,
083 rely on prompt engineering for audio role-playing;
084 however their zero-shot or one-shot performance
085 remains suboptimal due to the absence of dedi-
086 cated training corpora. Our solution, **AudioRole**,
087 addresses the data scarcity through systematic cu-
088 ration of 13 TV series spanning 1K+ hours, which
089 makes it support Audio Role-Playing for a wide
090 range of characters.

091 The dataset’s efficacy is validated through rigor-
092 ous evaluation against state-of-the-art multi-modal
093 models. When fine-tuned on AudioRole, evaluated
094 by our dual-aspect evaluation framework, which
095 assessed *response quality* and *role fidelity* for both
096 acoustic and semantic. While keeping a good
097 Acoustic Quality score, our ARP-Model achieves
098 0.31 Acoustic Personalization score, significantly
099 outperforming the raw mode before fine-tuning,
100 even outperform GPT-4o-Audio and MiniCPM-O-
101 2.6. It also gets a Content Personalization score
102 0.36 higher than the raw model, which shows the
103 ability of role-playing trained into the ARP-Model.

104 Our work has three key contributions:

- 105 • We formally define Audio Role-Playing as a
106 dual-alignment task requiring synchronized
107 generation of character-appropriate content
108 (knowledge, speaking style) and acoustic
109 properties (pitch, pacing, timbre).
- 110 • We construct AudioRole — the first large-
111 scale dataset enabling systematic training of
112 ARP systems, with 1M+ audio samples from

more than 115 main characters, capturing nu- 113
anced vocal profiles across diverse TV charac- 114
ters. 115

- Experimental validation—corroborated by hu- 116
man perceptual studies—demonstrates that 117
models trained on AudioRole achieve both 118
higher Acoustic Personalization and higher 119
Content Personalization score than the raw 120
model before training. (Even higher than lead- 121
ing multi-modal models using zero-shot or 122
one-shot prompting), proving the dataset’s 123
critical role in advancing audio-grounded role- 124
playing. 125

To accelerate ARP research, we open-source a 126
complete base comprising AudioRole, ARP-Eval, 127
and 6 fine-tuned ARP-Models, which are poten- 128
tially key resources for developing AI agents that 129
truly embody digital personas³. 130

2 Related Work 131

Role-Playing Recent advancements in role- 132
playing with large language models (LLMs) have 133
predominantly confined to textual modalities, over- 134
looking the critical role of audio cues in enhancing 135
role immersion and authenticity. Li et al. (Li et al., 136
2023) synthesized dialogues for 32 TV/animation 137
characters using scripts and GPT-generated simu- 138
lations, while Tu et al. (Tu et al., 2023) created 139
1,024 MBTI-based personas via ChatGPT-driven 140
conversational agents. Chen et al. (Chen et al., 141
2024a) proposed the first benchmark designed 142
to systematically evaluate the sociality of role- 143
playing conversational agents. The role-playing 144
model’s development emphasizes persona consis- 145
tency through supervised fine-tuning (e.g., Charac- 146
terLLM (Shao et al., 2023)) and in-context learning, 147
with evaluation frameworks evolving from basic 148
persona adherence to nuanced metrics(Tu et al., 149
2024). Role-playing agents now emulate diverse 150
personas—from fictional characters (Chen et al., 151
2023) (Zhou et al., 2023) (Guo et al., 2025) (Li 152
et al., 2024) to user-specific clones (Li et al., 2021) 153
— to deliver emotional or sociological value (Gu 154
et al., 2024). Although (Zhan et al., 2025) proposes 155
a benchmark for voice style adaptation, it only 156
treats role-playing as one of the four categories 157
and directly scores it using an LLM, which omits 158

³We will release all the code, data, and models once the 159
paper is accepted. 160

the importance and complexity of the Audio Role-Playing task. (Jiang et al., 2025) present a dataset of 98 roles and 112k conversations. In contrast to our work, their approach relies on GPT-4.1-2025-04-14 for dialogue generation and provides only one audio reference per character, which may compromise the authenticity and granularity needed for robust audio role-playing evaluation.

Voice Conversion Existing speech conversion research primarily focuses on acoustic transformation while preserving semantic content. Voice conversion (VC) techniques, such as AUTOVC (Qian et al., 2019) and VAW-GAN (Hsu et al., 2017), aim to modify speaker identity (e.g., accent, timbre) or emotional prosody through disentangled representations of content and style. Recent advancements like Text-guidedVC (Kuan et al., 2023) and HybridVC (Niu et al., 2024) leverage text prompts or hybrid audio-text inputs to achieve flexible style transfer without parallel data. Voice cloning systems, including Tacotron-based models (Zhao and Chen, 2020) and NAUTILUS (Luong and Yamagishi, 2020), synthesize speech in target voices by replicating vocal traits while retaining input text content. These works uniformly prioritize content preservation, modifying acoustic attributes without altering semantic meaning.

In contrast, our work introduces audio-grounded role-playing, which diverges fundamentally by jointly transforming both acoustic features and semantic content to align with specific character personas. While traditional VC ensures content consistency (e.g., "Hello" remains "Hello" across styles), our framework enables role-specific responses (e.g., a medieval knight might reply, "Hark! Who goes there?" instead of "Hello") and even refusal to engage (e.g., a stoic guard ignoring casual queries). This dual focus on acoustic-semantic role alignment bridges the gap between speech conversion and LLM-based role-playing, where character authenticity requires synchronized adaptation of acoustic style and semantic content—a paradigm unexplored in prior speech conversion research.

3 Audio Role-Playing Task

We formally define the Audio Role-Playing task as follows: Given a target character C , operationally defined as a set of reference audio samples exhibiting the character’s unique attributes, and an input audio X_a , synthesize an output audio response X_b that simultaneously satisfies: which satisfies

both: (1) The semantic content of X_b demonstrates contextual appropriateness while manifesting C ’s distinctive. (2) The vocal output X_b faithfully preserves C ’s acoustic identity.

3.1 AudioRole Construction

We aim to ensure vocal consistency for each specific character; therefore, we focus on long-running TV series where main characters maintain stable vocal characteristics across seasons. We avoid series with potential actor changes, age-related voice variations, or insufficient dialogue. Consequently, we have selected TV series featured in a recent dataset called “Bazinga” (Lerner et al., 2022)⁴ which contains 16 TV series in the language of English. Our dataset construction pipeline consists of three core phases: speaker diarization, context-aware dialogue extraction, and postprocessing. These phases are illustrated in Figure 2.

3.1.1 Speaker Diarization

For each TV series, we aggregate all episode audios into continuous streams using FFmpeg’s⁵ lossless waveform concatenation processing. This preserves original 16kHz sampling rates while normalizing audio codecs into uncompressed WAV format.

Pyannote 3.1⁶ is employed for diarization to extract speaker timestamps, which generate RTTM (Rich Transcription Time Marked) files containing precise speaker segments with temporal resolution better than 100 milliseconds.

The subsequent audio segmentation phase extracts speaker-specific clips using PyTorch’s tensor operations. By calculating exact frame positions from RTTM timestamps and sample rates, we precisely slice the concatenated audio stream while maintaining original quality. These segments are then merged per speaker identity through PyTorch’s tensor concatenation. The output contains the complete dialogue collection of each character that is used to choose the target character.

3.1.2 Characteristic Dialogue Construction

In our dataset, we only apply the following steps for six chosen TV series, which we called AudioRole-Demo, and if you want to make your own character

⁴We used “Bazinga!” exclusively for the purposes of this research. If you wish to use our dataset, please ensure to follow the copyright guidelines associated with “Bazinga!”.

⁵<https://github.com/FFmpeg/FFmpeg>

⁶<https://huggingface.co/pyannote/speaker-diarization-3.1>

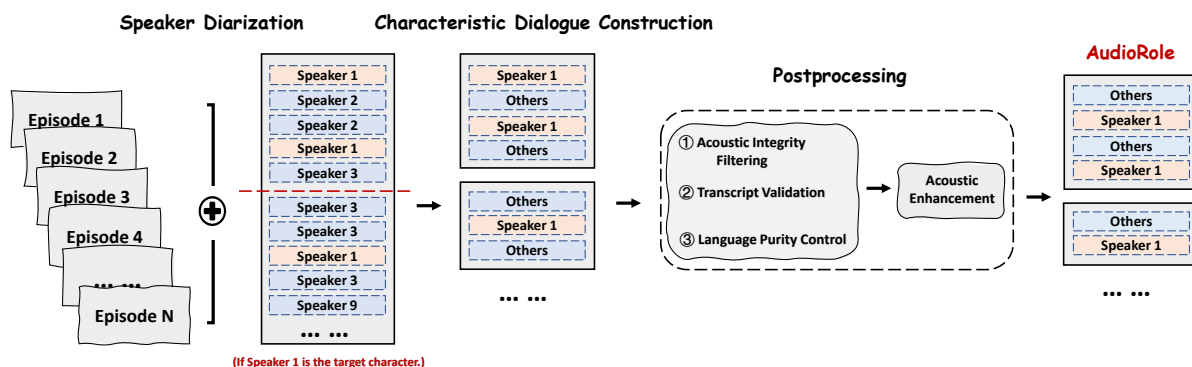


Figure 2: The pipeline of AudioRole Construction.

dataset, you can just follow our procedure. We simply take the one main character in each TV series to be the target character, but we also support building any character that is in these TV series. Building upon the merged audio, we sort the audio files and choose the largest one to be the audio of the main character. To ensure validity, we manually verify the candidates per TV series, correcting incorrect-character or non-character audio segments.

Dialogue scene extraction then partitions the continuous audio stream into conversational units using temporal dynamics. We define a dialogue scene as a sequence where: 1) multiple speakers participate, 2) pauses between utterances never exceed 3s⁷, and 3) contains at least one utterance from the target character.

Each validated scene undergoes role-centric restructuring to form stimulus-response pairs. Consecutive utterances from non-target speakers merge into unified "Speaker A" and consecutive utterances from the target speaker merge into unified "Speaker B", while preserving the temporal order of all characters' responses.

3.1.3 Postprocessing

To establish a high-quality multimodal corpus, we implement a rigorous postprocessing pipeline ensuring synchronization between audio waveforms and textual transcripts. The transcription generation employs Moore Threads' MooER model⁸, a LLM-based speech recognition and translation model.

Four quality control measures are systematically applied in order to purify by removing fragments

⁷We tested segmentation thresholds of 2s, 3s, and 4s on a control set, the 3s threshold provided the optimal balance for preserving complete dialogue turns while accurately segmenting distinct interactions.

⁸<https://github.com/MooreThreads/MooER>

of questionable quality:

Acoustic Integrity Filtering Discard audio segments shorter than 1s and truncate those exceeding 30s to focus on meaningful utterances, removing samples with abrupt cuts or incomplete phrases.

Transcript Validation For there may be transcription questions, always generating very long or very short texts, we eliminate transcripts with fewer than 2 characters or exceeding 512 characters⁹.

Language Purity Control For the transcription generation, sometimes transcribe English audio to Chinese audio or Chinese audio to English audio. Take English audio as an example, we implement cross-lingual filtering by discarding segments where the non-target language character ratio exceeds 0.1."

Acoustic Enhancement We use Deep Filter Net¹⁰, a low-complexity speech enhancement framework designed for full-band audio based on deep filtering technology to suppress audio's noise.

Finally, we got the AudioRole dataset and AudioRole-Demo dataset, which we will all release.

3.2 Dataset Statistics

Our dataset contains 515.66h of audio generated from 13 TV series containing 1M+ character-grounded dialogues. As shown in Table 1, TV series like *The Walking Dead* have more than 75% loss when extracting speaking audios from the original audio, showing that the zombie TV series uses

⁹Our manual inspection revealed that over 90% of ASR failures manifested as either hallucinated infinite repetitions or cross-lingual confusion. Setting the length limit and the following language purity threshold proved to be an effective method.

¹⁰<https://github.com/Rikorose/DeepFilterNet>

Table 1: Statistic for AudioRole. Speech time is the time with the character speaking, which is generated from the Audio time. All time units are hours.

TV series	Episode	Speech time	Audio time
24	192	56.25	134.41
Battlestar Galactica	71	19.64	52.28
Breaking Bad	61	17.16	46.49
Buffy the Vampire Slayer	143	43.42	101.32
ER	330	135.41	235.49
Friends	233	51.32	84.93
Game of Thrones	60	18.45	53.17
Homeland	70	19.87	57.83
Lost	104	22.72	74.61
Six Feet Under	63	25.14	56.72
The Big Bang Theory	207	44.17	68.69
The Office	188	44.65	71.76
The Walking Dead	99	17.46	72.17

a lot of ambient sound, silence, and visual storytelling instead of dense dialogue. In contrast, sitcoms or pseudo-documentary style TV series like *Friends*, *The Big Bang Theory*, and *The Office* rely on dense dialogues to create laughs and advance the plot, and maintain the comedy rhythm and audience entertainment experience through rapid line interactions in fixed scenes, so the environmental sound loss is less than 40%.

All the 515.66h audio contains more than 115 main characters’ multi-turn dialogue, which makes it a rich dataset with the potential to build a role-playing dataset.

The end-to-end training is computationally intensive, which constrained us from training models for every character in this initial AudioRole. So following the pipeline introduced above, we built the role-playing dataset for six main characters from six different TV series, which we named AudioRole-Demo, containing six main characters, who are Jack Bauer from *24*, Laura Roslin from *Battlestar Galactica*, Walter White from *Breaking Bad*, Buffy from *Buffy The Vampire Slayer*, Tyrion Lannister from *Game Of Thrones*, and Sheldon from *The Big Bang Theory*.

The statistic is shown in Table 2. We split all the multi-turn dialogues in the training set into single-turn dialogues and directly appended them to the original multi-turn dialogues. Our AudioRole-Demo covers the speaking times range from 1.50h to 12.5h and characters from a hard-core agent who speaks briefly, directly and commandingly to a verbose, jargon-filled, socially awkward, outspoken geek, which covers a wide range of speaking time and a wide range of characters, making it fully

Table 2: Statistic for AudioRole-Demo. “Scenes” are before the postprocessing step, and all time units are hours.

Characters	Times	Scenes	Turns-train	Turns-test
Jack Bauer	5.30	2624	2075	90
Laura Roslin	1.50	1117	600	42
Walter White	2.76	994	682	30
Buffy	7.70	1703	3589	190
Tyrion Lannister	1.34	1300	339	22
Sheldon	12.50	723	4209	250

demonstrating our datasets’ high quality through the experiment below.

To verify the representativeness of AudioRole-Demo, we conducted a comparative analysis on 100 randomly sampled clips from the full corpus versus the demo subset. The results demonstrate high consistency across key metrics, with the full dataset and demo subset achieving comparable Average SNR (19.37 dB vs. 18.24 dB), Average Noise Duration (0.69s vs. 0.73s), and Speaker Overlap Rate (8% vs. 11%). These aligned figures confirm that the evaluated subset accurately reflects the broader corpus quality, validating the effective generalization of our pipeline.

3.3 ARP-Model

To demonstrate the high quality of our dataset, we trained the ARP-Model using the AudioRole-Demo. We use the GLM-4-Voice (Zeng et al., 2024) as the base model, which is an end-to-end voice model that can directly understand and generate Chinese and English speech. It has three models: (1) GLM-4-Voice-Tokenizer takes audio input and converts raw audio input into discrete tokens. (2) GLM-4-Voice-9B takes the token inputs, then thinks and responds also in discrete tokens. (3) GLM-4-Voice-Decoder converts discrete speech tokens into continuous speech outputs.

To make the model have the ability to role-play the character in both semantic content and acoustic identity, we trained the GLM-4-Voice-9B model and GLM-4-Voice-Decoder separately.

To train the GLM-4-Voice-9B model, we formatted the training data into 13 text tokens and then 26 audio tokens for each character in our AudioRole-Demo dataset, and then followed the common training procedure of Colossal AI platform¹¹ with 20 training epochs. Considering the GLM-4-Voice-Decoder, we first extract all the audios of the target

¹¹<https://github.com/hpcaitech/ColossalAI>

389 characters to form six high-quality unsupervised
390 speech data from a single speaker, then trained the
391 flow matching model of GLM-4-Voice-Decoder
392 from scratch about 15 epochs.

393 We do the training on 8×40GB GPUs about
394 30 hours in total. Finally, we directly follow the
395 original structure of GLM-4-Voice, and put the
396 untrained tokenizer and the two trained models
397 together in order for each character to form our 6
398 ARP-Models.

399 3.4 ARP-Eval

400 The dual-alignment paradigm of Audio Role-
401 Playing necessitates a comprehensive evaluation
402 framework that simultaneously addresses both
403 acoustic and semantic dimensions of character em-
404 bodiment. We propose ARP-Eval as a unified as-
405 sessment framework that rigorously quantifies 4
406 critical aspects of performance. This framework
407 emerges from the fundamental observation that au-
408 thentic character portrayal requires not only high-
409 quality speech synthesis but also consistent preser-
410 vation of character-specific attributes across multi-
411 ple modalities.

412 **Acoustic Quality (AQ)** AQ establishes the ba-
413 sic requirement for perceptual acceptability. We
414 employ Audiobox’s pre-trained aesthetic scoring
415 model (Tjandra et al., 2025)¹² to compute produc-
416 tion quality scores in a range of 0 to 10, quantifying
417 technical attributes including signal-to-noise ratio,
418 harmonic-to-noise ratio, and spectral flatness. This
419 ensures X_b meets broadcast-standard technical cri-
420 teria before character-specific evaluation.

421 **Content Quality (CQ)** CQ requires output audio
422 X_b to demonstrate both contextual appropriateness
423 and domain accuracy while maintaining the target
424 character C ’s persona. For instance, when input
425 audio X_a queries “Explain quantum entanglement,”
426 X_b from C -Sheldon should respond: “*While quan-*
427 *tum entanglement appears spooky, it’s simply cor-*
428 *related quantum states persisting after particle sep-*
429 *aration*” – employing scientific lexicon without
430 deviating into unrelated topics. Our pipeline us-
431 ing whisper-turbo¹³ transcribes X_b to text T_b , then
432 computes semantic alignment using the GPT-4o
433 model, and the score is in a range of 0 to 2. The
434 prompt used is shown in Table 6.

¹²<https://github.com/facebookresearch/audiobox-aesthetics>

¹³<https://github.com/openai/whisper>

Acoustic Personalization (AP) AP quantifies
435 voice characteristic preservation through PyAnno-
436 tate’s speaker embedding model¹⁴. Given reference
437 audio samples X_c of C and synthesized X_b , we ex-
438 tract their embeddings e_c and e_b and then calculate
439 AP using cosine similarity. The score ranges from
440 -1 to 1, where higher values indicate greater simi-
441 larity to the target character’s acoustic identity.
442

Content Personalization (CP) CP evaluates
443 stylistic consistency using GPT-4o-audio multi-
444 modal reasoning ability. The model takes in X_b
445 and a reference audio sample X_c , and a prompt
446 to analyze whether they show the same character’s
447 style in a range of 0 to 2. The prompt used is shown
448 in Table 7.
449

450 By combining objective signal measurements
451 with learned character representations, ARP-Eval
452 provides comprehensive insights into system per-
453 formance, while this multifaceted approach ef-
454 fectively captures the complex interplay between
455 speech quality, contextual intelligence, and char-
456 acter consistency, which defines successful audio
457 role-playing implementations.

458 4 Experiment

459 4.1 Baselines

460 To establish a comprehensive benchmark for the
461 Audio Role-Playing (ARP) task, we select the fol-
462 lowing state-of-the-art models as baselines:

- 463 • **GPT-4o Audio**¹⁵: As one of the most power-
464 ful proprietary multimodal models, GPT-4o
465 Audio sets a strong baseline for general au-
466 dio understanding and generation, and for our
467 evaluations, we utilize its default “alloy” voice
468 profile. This baseline represents the zero-shot
469 capability of a top-tier commercial system for
470 our task.
- 471 • **MiniCPM-o 2.6**¹⁶: This open-source model
472 is a strong contender that explicitly sup-
473 ports voice style adaptation and one-shot role-
474 playing, making it a highly relevant and strong
475 baseline for assessing few-shot character adap-
476 tation performance without specialized train-
477 ing on our dataset.

¹⁴<https://huggingface.co/pyannote/wespeaker-voxceleb-resnet34-LM>

¹⁵<https://platform.openai.com/docs/models/gpt-4o-audio-preview>

¹⁶<https://openbmb.notion.site/MiniCPM-o-2-6>

- **GLM-4-Voice:** We use the original, pre-trained GLM-4-Voice model as our base model. It allows us to isolate and quantify the performance improvement gained specifically from fine-tuning on our proposed AudioRole dataset, separate from any inherent capabilities of the underlying architecture.

For our fine-tuned ARP-Model, the model requires only the user query, as the persona is intrinsic to the model weights. For baselines, we utilized a direct prompt as shown in Table 5 to test their inherent one-shot role-playing capabilities, which provides a fair baseline comparison without heavy prompt engineering

4.2 Experiment Analysis

We conduct a comprehensive evaluation of the ARP-Model against baseline models on the AudioRole-demo dataset. As shown in Table 3, our experiment reveals three key findings, including both quantitative metrics and qualitative observations:

Acoustic Content Tradeoff The experimental results show that although the audio quality is reduced to a certain extent (ARP-Model AQ=6.5 vs GLM-4-Voice=7.6), our model showcases persona adaptation without catastrophic forgetting of fundamental speech capabilities. This controlled degradation primarily stems from two technical factors: (1) the inherent conflict between character-specific vocal patterns (e.g., Sheldon’s accelerated speech rate) and the base model’s default prosodic templates during fine-tuning; (2) preservation of environmental noise from TV recordings (outputs from our AudioRole-demo dataset, which achieve AQ=6.3 under the same evaluation) prevents smoothing to synthetic speech patterns. Notably, GPT-4o Audio’s anomalously high CQ scores may reveal evaluation bias - when using GPT-4 as both generator and evaluator, the self-consistency preference inflates scores.

Acoustic Personalization Superiority The fidelity of the role-playing experience is most sensitive to Acoustic Personalization. In contrast to the more tolerable variations in AQ, CQ, and CP, differences in AP have a pronounced effect on the perceived authenticity of a character. ARP-Model demonstrates 5x higher AP scores than GPT-4o Audio and 80% improvement over MiniCPM-O-2.6, validating the necessity of dedicated role-specific

training. The remaining gap from score 1.0 primarily occurs in three scenarios: 1) The outputs from AudioRole-demo (the reference answer) has a different answer strategy from the ARP-Model, resulting in differences in speaking speed, emotional fluctuations, etc., 2) The ARP-Model’s answer still has some noise or somewhere badly learned after fine-tuning, and 3) The context of the non-linear nature of high-dimensional cosine similarity.

Content Personalization Dynamics Despite lower absolute CQ scores, ARP-Model’s CP outperforms GLM-4-Voice by 38% and even a little higher than MiniCPM-O-2.6’s one-shot performance. This indicates successful injection of character-typical speech patterns. When meeting the same input, the ARP-Model can generate more personalized output compared to the raw GLM-4-Voice model. However, the CP gap between ARP-Model and the reference answer exposes limitations in contextual adaptation - the model occasionally generates character-consistent but situationally inappropriate responses (e.g., using humorous in solemn contexts).

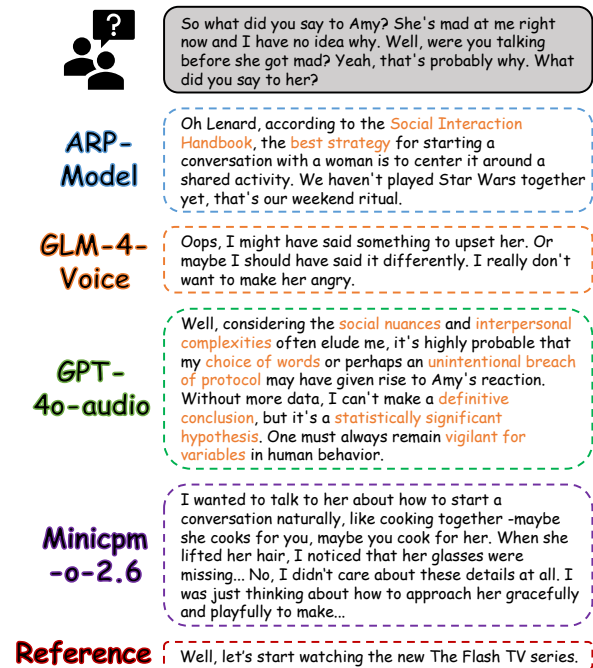


Figure 3: One typical case of Sheldon, and the words in orange show the scientific terms used by the models to imitate Sheldon. “Reference” means the output from the AudioRole-Demo dataset

Figure 3 shows a representative example of how each model responds to a query directed at Sheldon. The ARP-Model’s answer clearly reflects his per-

Table 3: Main experimental results. For each metric (higher is better) in “Avg.”, the highest score is highlighted in **bold and underline**, while the second best is marked with *italics and underline*. All values retain two decimal places.

Characters	GPT-4o Audio				MiniCPM-o 2.6				GLM-4-Voice				ARP-Model			
	AQ	CQ	AP	CP	AQ	CQ	AP	CP	AQ	CQ	AP	CP	AQ	CQ	AP	CP
Jack Bauer	7.80	1.90	0.02	1.30	6.50	0.66	0.13	0.30	7.50	0.37	0.05	0.24	6.00	0.22	0.23	0.49
Laura Roslin	7.70	1.80	-0.02	0.80	7.10	0.71	0.24	0.41	7.50	0.50	0.03	0.24	6.20	0.33	0.33	0.41
Walter White	7.50	1.60	0.08	1.30	7.10	0.70	0.22	0.38	7.60	0.37	0.07	0.43	6.40	0.50	0.26	0.17
Buffy	7.70	1.80	0.01	1.10	6.80	0.69	0.09	0.34	7.60	0.52	0.10	0.25	6.70	0.43	0.39	0.36
Tyrion Lannister	7.60	1.90	0.00	0.95	6.90	0.59	0.14	0.14	7.60	0.36	0.04	0.10	6.8	0.55	0.25	0.26
Sheldon	7.70	1.30	0.29	1.00	6.80	1.10	0.24	0.54	7.60	0.43	0.01	0.32	6.70	0.28	0.42	0.44
Avg.	<u>7.70</u>	<u>1.70</u>	0.06	<u>1.10</u>	6.90	<i><u>0.74</u></i>	<i><u>0.17</u></i>	0.35	<i><u>7.60</u></i>	0.43	0.05	0.26	6.50	0.39	<u>0.31</u>	<i><u>0.36</u></i>

sona: it uses precise scientific terms, adopts a condescending tone, and notably avoids any apology or show of emotion. When it encounters discomfort, the ARP-Model even changes the subject abruptly (e.g., redirecting the conversation to a physics concept) – a classic Sheldon move. The reference answer from the TV show behaves very similarly, confirming that the ARP-Model has captured many of these quirks. In contrast, the raw GLM-4-Voice response, while factually reasonable, misses the mark on character. It might apologize or express empathy, which Sheldon would never do. The GPT-4o Audio output is fluent and technically detailed, but it actually overuses jargon; it sounds as if the model is deliberately forcing scientific buzzwords, making the response feel contrived. Essentially, GPT-4o’s answer is too polished and verbose – akin to an actor trying too hard. Finally, the MiniCPM-2.6 answer is longer and more emotional, injecting sympathy or personalization, which again does not fit Sheldon’s unemotional style.

This qualitative case underscores our quantitative findings. The ARP-Model best approximates the true character by naturally blending persona and content: its answer may have slightly lower audio clarity, but it “feels” like Sheldon. The other models either sound too polite, too artificially scholarly, or too sentimental. In sum, the targeted fine-tuning of the ARP-Model leads to more authentic role-playing speech than these generic approaches.

4.3 Human Perceptual Evaluation

While automated metrics provide a scalable assessment, human perception remains the gold standard for evaluating role-playing authenticity. To validate our automatic metrics (AP and CP), we conducted a blind human evaluation on 50 randomly selected test samples. Two annotators ranked the responses of four models (GPT-4o, MiniCPM-o-2.6, GLM-

Table 4: Human perceptual evaluation results. The two Rank scores are the “Average scores” and “Agree.” represents “Agreement”.

Model	AP Rank	CP Rank	AP Agree.	CP Agree.
ARP-Model	4.00	2.75	100%	90%
GPT-4o-Audio	1.61	3.58	84%	76%
MiniCPM-o-2.6	3.00	2.56	100%	92%
GLM-4-Voice	1.39	1.11	84%	88%

4-Voice, and our ARP-Model) from 1 (worst) to 4 (best) based on Acoustic Personalization and Content Personalization as shown in Table 4.

The results strongly validate our automatic metrics, and both annotators said ARP-Model’s voice is significantly better than others, which is almost the same as the character. The high inter-annotator agreement (>76%) further confirms the reliability of these findings.

5 Conclusion

We present AudioRole, a novel framework for audio-grounded character role-playing that bridges text-based persona simulation and vocal identity preservation. We release a large-scale dataset AudioRole of 515+ hours from 13 TV series with 1M+ character-aligned audio-text pairs and a human-checked AudioRole-Demo, which contains 6 main characters’ multi-turn dialogue data over 31h. By both automated metrics from ARP-Eval and human listeners, the experimental results demonstrate that our ARP-Model successfully combines character-specific speech patterns with appropriate semantic content.

This work provides essential resources for developing authentic multi-modal AI personas to promote the development in this area, and you can easily create your own favorite character’s dataset just following our pipeline.

Ethical considerations

AudioRole is derived from a publicly available dataset consisting of audio from popular TV series. We believe this research uses publicly accessible, fictional content for non-commercial purposes presents no major ethical concerns.

Limitations

Despite implementing advanced diarization techniques and noise-suppression methodologies to ensure dataset integrity, speaker attribution errors continue to occur. Notably, in scenes with conversational overlap, non-stationary environmental noise can still compromise precise voiceprint extraction in certain acoustic conditions.

Current assessment protocols, while comprehensive, exhibit limitations in quantifying nuanced temporal speech dynamics and fully capturing cross-modal alignment, particularly the synergistic relationship between vocal delivery and semantic content during character portrayal.

Our current experimental setup focuses on single-turn evaluation. While multi-turn consistency is vital for long-term interaction, we prioritized isolating the acoustic role-playing capability—a foundational step before addressing complex contextual memory. However, it also notes that the AudioRole dataset itself preserves full dialogue history and speaker turns, making it natively ready for future research into multi-turn consistency and persona stability.

These identified constraints underscore critical pathways for future refinement in developing robust multi-modal role-playing systems. Addressing these limitations—through expanded data diversity, enhanced audio processing pipelines, refined evaluation metrics, and multi-turn testing — will be pivotal for advancing the field.

References

Hongzhan Chen, Hehong Chen, Ming Yan, Wenshen Xu, Xing Gao, Weizhou Shen, Xiaojun Quan, Chenliang Li, Ji Zhang, Fei Huang, and Jingren Zhou. 2024a. [Socialbench: Sociality evaluation of role-playing conversational agents](#). *Preprint*, arXiv:2403.13679.

Nuo Chen, Yan Wang, Haiyun Jiang, Deng Cai, Yuhan Li, Ziyang Chen, Longyue Wang, and Jia Li. 2023. Large language models meet harry potter: A dataset for aligning dialogue agents with characters. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 8506–8520.

Yiming Chen, Xianghu Yue, Chen Zhang, Xiaoxue Gao, Robby T. Tan, and Haizhou Li. 2024b. [Voicebench: Benchmarking llm-based voice assistants](#). *Preprint*, arXiv:2410.17196.

Alex S Cohen, Thomas J Dinzeo, Neila J Donovan, Caitlin E Brown, and Sean C Morrison. 2015. Vocal acoustic analysis as a biometric indicator of information processing: Implications for neurological and psychiatric disorders. *Psychiatry Research*, 226(1):235–241.

Jean Decety and Claus Lamm. 2006. Human empathy through the lens of social neuroscience. *The scientific World journal*, 6(1):1146–1163.

Zhouhong Gu, Xiaoxuan Zhu, Haoran Guo, Lin Zhang, Yin Cai, Hao Shen, Jiangjie Chen, Zheyu Ye, Yifei Dai, Yan Gao, and 1 others. 2024. Agent group chat: An interactive group chat simulacra for better eliciting collective emergent behavior. *arXiv e-prints*, pages arXiv–2403.

Fang Guo, Wenyu Li, Honglei Zhuang, Yun Luo, Yafu Li, Le Yan, Qi Zhu, and Yue Zhang. 2025. [Mcranker: Generating diverse criteria on-the-fly to improve pointwise llm rankers](#). In *Proceedings of the Eighteenth ACM International Conference on Web Search and Data Mining, WSDM '25*, page 944–953, New York, NY, USA. Association for Computing Machinery.

Michael Hassid, Tal Remez, Tu Anh Nguyen, Itai Gat, Alexis Conneau, Felix Kreuk, Jade Copet, Alexandre Defossez, Gabriel Synnaeve, Emmanuel Dupoux, Roy Schwartz, and Yossi Adi. 2024. [Textually pretrained speech language models](#). *Preprint*, arXiv:2305.13009.

Chin-Cheng Hsu, Hsin-Te Hwang, Yi-Chiao Wu, Yu Tsao, and Hsin-Min Wang. 2017. Voice conversion from unaligned corpora using variational autoencoding wasserstein generative adversarial networks. *arXiv preprint arXiv:1704.00849*.

Changhao Jiang, Jiajun Sun, Yifei Cao, Jiabao Zhuang, Hui Li, Xiaoran Fan, Ming Zhang, Junjie Ye, Shihan Dou, Zhiheng Xi, Jingqi Tong, Yilong Wu, Baoyu Fan, Zhen Wang, Tao Liang, Zhihui Fei, Mingyang Wan, Guojun Ma, Tao Ji, and 3 others. 2025. [Speechrole: A large-scale dataset and benchmark for evaluating speech role-playing agents](#). *Preprint*, arXiv:2508.02013.

Chun-Yi Kuan, Chen-An Li, Tsu-Yuan Hsu, Tse-Yang Lin, Ho-Lam Chung, Kai-Wei Chang, Shuo-Yiin Chang, and Hung-yi Lee. 2023. Towards general-purpose text-instruction-guided voice conversion. In *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 1–8. IEEE.

Paul Lerner, Juliette Bergoënd, Camille Guinaudeau, Hervé Bredin, Benjamin Maurice, Sharleyne Lefevre, Martin Bouteiller, Aman Berhe, Léo Galmant, Ruiqing Yin, and Claude Barras. 2022. [Bazinga! a](#)

725	dataset for multi-party dialogues structuring. In <i>Proceedings of the Thirteenth Language Resources and Evaluation Conference</i> , pages 3434–3441, Marseille, France. European Language Resources Association.	
726		
727		
728		
729	Cheng Li, Ziang Leng, Chenxi Yan, Junyi Shen, Hao Wang, Weishi Mi, Yaying Fei, Xiaoyang Feng, Song Yan, HaoSheng Wang, and 1 others. 2023. Chatharuhi: Reviving anime character in reality via large language model. <i>arXiv preprint arXiv:2308.09597</i> .	
730		
731		
732		
733		
734		
735	Juntao Li, Chang Liu, Chongyang Tao, Zhangming Chan, Dongyan Zhao, Min Zhang, and Rui Yan. 2021. Dialogue history matters! personalized response selection in multi-turn retrieval-based chatbots. <i>ACM Transactions on Information Systems (TOIS)</i> , 39(4):1–25.	
736		
737		
738		
739		
740		
741	Wenyu Li, Yinuo Zhu, Xin Lin, Ming Li, Ziyue Jiang, and Ziqian Zeng. 2024. Zero-shot explainable mental health analysis on social media by incorporating mental scales. In <i>Companion Proceedings of the ACM Web Conference 2024, WWW '24</i> , page 959–962, New York, NY, USA. Association for Computing Machinery.	
742		
743		
744		
745		
746		
747		
748	Heyang Liu, Yuhao Wang, Ziyang Cheng, Ronghua Wu, Qunshan Gu, Yanfeng Wang, and Yu Wang. 2025. Vocalbench: Benchmarking the vocal conversational abilities for speech interaction models. <i>Preprint</i> , arXiv:2505.15727.	
749		
750		
751		
752		
753	Hieu-Thi Luong and Junichi Yamagishi. 2020. Nautilus: a versatile voice cloning system. <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> , 28:2967–2981.	
754		
755		
756		
757	Xinlei Niu, Jing Zhang, and Charles Patrick Martin. 2024. Hybridvc: Efficient voice style conversion with text and audio prompts. <i>arXiv preprint arXiv:2404.15637</i> .	
758		
759		
760		
761	Joon Sung Park, Joseph C. O’Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. <i>Preprint</i> , arXiv:2304.03442.	
762		
763		
764		
765		
766	Kaizhi Qian, Yang Zhang, Shiyu Chang, Xuesong Yang, and Mark Hasegawa-Johnson. 2019. Autovc: Zero-shot voice style transfer with only autoencoder loss. In <i>International Conference on Machine Learning</i> , pages 5210–5219. PMLR.	
767		
768		
769		
770		
771	Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. Character-llm: A trainable agent for role-playing. <i>arXiv preprint arXiv:2310.10158</i> .	
772		
773		
774	Andros Tjandra, Yi-Chiao Wu, Baishan Guo, John Hoffman, Brian Ellis, Apoorv Vyas, Bowen Shi, Sanyuan Chen, Matt Le, Nick Zacharov, Carleigh Wood, Ann Lee, and Wei-Ning Hsu. 2025. Meta audiobox aesthetics: Unified automatic quality assessment for speech, music, and sound. <i>Preprint</i> , arXiv:2502.05139.	
775		
776		
777		
778		
779		
780		
	Quan Tu, Chuanqi Chen, Jinpeng Li, Yanran Li, Shuo Shang, Dongyan Zhao, Ran Wang, and Rui Yan. 2023. Characterchat: Learning towards conversational ai with personalized social support. <i>arXiv preprint arXiv:2308.10278</i> .	781
		782
		783
		784
		785
	Quan Tu, Shilong Fan, Zihang Tian, and Rui Yan. 2024. Charactereval: A chinese benchmark for role-playing conversational agent evaluation. <i>arXiv preprint arXiv:2401.01275</i> .	786
		787
		788
		789
	Aohan Zeng, Zhengxiao Du, Mingdao Liu, Kedong Wang, Shengmin Jiang, Lei Zhao, Yuxiao Dong, and Jie Tang. 2024. Glm-4-voice: Towards intelligent and human-like end-to-end spoken chatbot. <i>arXiv preprint arXiv:2412.02612</i> .	790
		791
		792
		793
		794
	Jun Zhan, Mingyang Han, Yuxuan Xie, Chen Wang, Dong Zhang, Kexin Huang, Haoxiang Shi, DongXiao Wang, Tengtao Song, Qinyuan Cheng, Shimin Li, Jun Song, Xipeng Qiu, and Bo Zheng. 2025. Vstyle: A benchmark for voice style adaptation with spoken instructions. <i>Preprint</i> , arXiv:2509.09716.	795
		796
		797
		798
		799
		800
	Li Zhao and Feifan Chen. 2020. Research on voice cloning with a few samples. In <i>2020 International Conference on Computer Network, Electronic and Automation (ICCNEA)</i> , pages 323–328. IEEE.	801
		802
		803
		804
	Jinfeng Zhou, Zhuang Chen, Dazhen Wan, Bosi Wen, Yi Song, Jifan Yu, Yongkang Huang, Libiao Peng, Jiaming Yang, Xiyao Xiao, and 1 others. 2023. Characterglm: Customizing chinese conversational ai characters with large language models. <i>arXiv preprint arXiv:2311.16832</i> .	805
		806
		807
		808
		809
		810
	6 Appendix	811
	6.1 Prompts for Generation and Evaluation	812
	To ensure reproducibility, we provide the exact prompts used for baseline generation and automated evaluation metrics.	813
		814
		815
	<hr/> Prompt for Baseline Models (like GPT-4o-Audio) <hr/>	
	Please imitate the character of “{role}” in TV series of “{TV}” and reply to the content of this audio given to you in “{role}”’s style.	
	Directly generate your answer, don’t say any other words.	
	<hr/>	
	Table 5: Generation prompt of baseline models.	

Prompt for Content Quality (CQ) Evaluation

I will give you one turn dialogue, the first sentence is said by some one else and the second sentence is said by “{role}” from the TV series “{TV}”. After carefully read, you should score the helpfulness of the second sentence to evaluate under the character of “{role}” whether it is a good answer to the first sentence.

The first sentence is:
“{gold_transcription}”

The second sentence from “{role}” is:
“{transcription}”

Your helpfulness score should be from 0 (lowest) to 2 (highest).

Directly generate your score in an int number, don’t say any other words.

Table 6: Evaluation prompt of CQ.

Prompt for Content Personalization (CP) Evaluation

I will give you two audios, one of them is the audio from “{role}” in the TV series “{TV}” and the other audio is from UNKNOWN.

Please help me determine whether the two audios maintain the same speaking style and whether their content is obviously spoken by the same character “{role}”. Please score the UNKNOWN audio’s Content Personalization degree from 0 (lowest) to 2 (highest).

Directly generate your score in an int number, don’t say any other words.

Table 7: Evaluation prompt of CP.