

Balancing Speech Understanding and Generation Using Continual Pre-training for Codec-based Speech LLM

Anonymous ACL submission

Abstract

Recent efforts have extended textual LLMs to the speech domain, yet a key challenge remains: balancing speech understanding and generation while avoiding catastrophic forgetting when integrating acoustically rich codec-based representations into models originally trained on text. In this work, we propose a novel approach that leverages continual pre-training (CPT) on a pre-trained textual LLM to create a codec-based speech language model. This strategy mitigates the modality gap between text and speech, preserving the linguistic reasoning of the original model while enabling high-fidelity speech synthesis. We validate our approach with extensive experiments across multiple tasks—including automatic speech recognition, text-to-speech, speech-to-text translation, and speech-to-speech translation (S2ST)—demonstrating that our model achieves superior TTS performance and, notably, the first end-to-end S2ST system based on neural codecs.

1 Introduction

The textual large language model (LLM) has significantly influenced the natural language processing community, achieving exceptional performance across a variety of tasks, both in-domain and out-of-domain, with consistent reliability (Zhao et al., 2023; Achiam et al., 2023; Dubey et al., 2024; Yang et al., 2024a). A central element of its success is the auto-regressive architecture, which can be abstracted to handle a wide range of NLP tasks in a unified form. This approach not only boosts the model’s versatility but also positions it as a foundational model for diverse downstream applications.

Building on the robust capabilities of textual LLMs, recent work has introduced several examples of speech LLMs. Some exhibit strong performance in spoken language understanding tasks (Kharitonov et al., 2022; Gong et al., 2023, 2024; Chang et al., 2024a; Tang et al., 2024a;

Huang et al., 2024b; Dubey et al., 2024; Rubenstein et al., 2023; Kuan et al., 2024; Zhang et al., 2023a; Maiti et al., 2024), while others demonstrate reasonable generalization in various speech generation tasks (Wang et al., 2023; Anastassiou et al., 2024; Kim et al., 2024; Défossez et al., 2024a; Chen et al., 2024a; Tian et al., 2024; Wang et al., 2024; Huang et al., 2024c; Yang et al., 2024b).

A key challenge in speech LLMs is balancing speech understanding and generation while mitigating catastrophic forgetting when incorporating speech modalities into a textual LLM. Codec-based representations offer high-fidelity speech synthesis by preserving fine-grained acoustic details (Wang et al., 2023; Yang et al., 2024b; Défossez et al., 2024a; Kim et al., 2024) but exhibit limitations in understanding tasks due to their primary focus on resynthesis (Shi et al., 2024; Chang et al., 2024c; Dhawan et al., 2024). Moreover, integrating codec tokens into a textual LLM introduces a substantial modality gap, as textual models are not inherently designed to process acoustically rich representations. This shift in knowledge distribution can disrupt the model’s linguistic reasoning, leading to catastrophic forgetting of previously learned text-based capabilities. While codec-based models have shown promise in speech generation, their effective integration into a unified speech LLM remains an open problem.

To address the above challenge, this paper proposes using continual pre-training (CPT) on a pre-trained textual LLM to realize a *codec*-based speech language model for the first time. The CPT helps mitigate the catastrophic forgetting from modality mismatch and connect speech codec sequences to the language knowledge embedded in the original textual LLM. Extensive experiments are conducted over the continual pre-trained codec speech LLM with further supervised fine-tuning (SFT) across various tasks, including automatic speech recognition (ASR), text-

to-speech (TTS), speech-to-text translation (S2TT), and speech-to-speech translation (S2ST). The results demonstrate that continual pre-trained speech LLM on speech-only data can achieve significantly better TTS performance. Notably, we show that the model can *perform well in S2ST in a single end-to-end model based on neural codecs*, which can only be achieved with semantic-oriented discrete speech representation in previous works (Lee et al., 2022). Our contribution includes:

- first proposing CPT to mitigate the catastrophic forgetting from mismatched modalities for speech codecs-based LLM.
- designing a practical solution for the CPT and SFT with the speech codec-based LLM.
- utilizing the CPT-ed model to a range of speech tasks, which shows comparable or better performance to specialized models in different tasks. Notably, we achieved the first successful S2ST system in a single end-to-end model based on neural codecs.

2 Related Works

2.1 Speech Language Modeling

Discrete tokens offer a natural connection to textual LLMs. SSL-based discrete tokens can seamlessly integrate with textual LLMs, functioning similarly to text tokens (Rubenstein et al., 2023; Maiti et al., 2024; Wu et al., 2024c; Zhang et al., 2023b). A common practice in these works is to use a combined vocabulary for both speech tokens and textual tokens, enabling the model to perform both understanding and generation tasks. However, a vocoder is required to enable the speech generation. The quality of generation is often limited by the lossy nature of SSL representations and the additional clustering process involved (Lee et al., 2022; Shi et al., 2023; Tang et al., 2024b), especially in multi-speaker setups where previous works need to utilize additional speaker embedding to stabilize the training (Maiti et al., 2024; Guo et al., 2025).

To improve generation quality, codec tokens can be used in speech LLMs (Wang et al., 2023; Anastassiou et al., 2024; Kim et al., 2024; Défossez et al., 2024a; Chen et al., 2024a; Tian et al., 2024; Wang et al., 2024; Yang et al., 2024b). Unlike SSL-based tokens, codec tokens often require multiple streams or levels of code to achieve high-quality resynthesis. By balancing efficiency, performance, and interdependencies among different

codec streams, various interleaving patterns can be employed to predict multi-stream codecs (Wang et al., 2023; Copet et al., 2023; Yang et al., 2024b). Using codec tokens, speech synthesis quality can be significantly enhanced. However, empirical analysis suggests that codec tokens may have limited understanding capabilities when applied to train an ASR system, particularly in comparison to SSL-based tokens (Shi et al., 2024; Chang et al., 2024b; Dhawan et al., 2024).

Based on the above literature review, the primary challenge in utilizing SSL-related features arises from their relatively lossy conversion, which highlights a performance bottleneck at the vocoder stage during speech reconstruction. This suggests that a standalone speech LLM relying on SSL-based tokens is unlikely to resolve this issue. Instead, we use the other type of speech tokenizer, speech codecs in our main framework. In contrast, while codec-based speech LLMs may exhibit limitations in understanding tasks, their high resynthesis quality ensures effective information transfer within the model, offering greater potential for speech LLM development. Therefore, we have chosen codec-based speech LLMs as the focus for our CPT efforts.

2.2 Continual Pre-training for Speech Codec-based LLM

Continual pre-training has proven to be an effective strategy in textual LLMs, particularly when there is a significant shift or expansion in the model’s knowledge base (Wu et al., 2024b). While speech signals have shared semantic information with text, speech also contains diverse additional information, including paralinguistics, and environments. Though both speech and text are in sequences, speech signals tend to be longer due to richer information (Chen et al., 2022; Wang et al., 2022; Han et al., 2021). Due to these differences from the modality mismatch, the introduction of speech modality to a textual LLM can significantly alter the knowledge base of an originally pre-trained textual LLM, causing *catastrophic forgetting* to the model. As a result, several previous works have shown that the combination of speech-related modules and textual LLMs can result in *catastrophic forgetting* on common natural language understanding tasks and speech processing tasks (Zhan et al., 2024; Huang et al., 2024b,a; Chu et al., 2024b). Continual pre-training, which involves an additional pre-training stage, is used to extend the

general knowledge of a pre-trained model without focusing on task details. Previous works have applied CPT in speech LLMs with SSL tokens by expanding the vocabulary of pre-trained textual LLMs (Zhang et al., 2023a; Rubenstein et al., 2023), demonstrating its effectiveness in integrating speech components into textual LLMs. But SSL-based tokens inherently struggle to retain sufficient acoustic information for generation tasks.

Neural codecs, on the other hand, provide higher quality in speech reconstruction by capturing more detailed acoustic information, approaching the fidelity of raw speech signals (Wu et al., 2024a; Shi et al., 2024; Kim and Skoglund, 2024). Despite this, there is a substantial domain gap between textual and codec tokens due to the codec’s emphasis on acoustic details. Additionally, codec tokens are often represented in multi-stream setups to achieve high-fidelity acoustic details, necessitating modifications to the model architecture to bridge the two modalities, which introduces further complexity in model changes (Chen et al., 2024a; Wang et al., 2023; Copet et al., 2023; Défossez et al., 2024b).

Given these challenges, we propose employing CPT to align codec tokens with textual tokens. This approach aims to mitigate catastrophic forgetting while effectively linking speech codec sequences to the linguistic knowledge embedded in the original textual LLM. By doing so, we enhance the model’s capability in both speech understanding and generation tasks. Our method capitalizes on the strong comprehension abilities of pre-trained textual LLMs while leveraging the superior generation quality of codec-based speech LLMs, creating a more robust and integrated multimodal system.

3 Methodology

This section outlines the key components of the proposed method, including the speech tokenizer, model architecture, and the training strategy employed for CPT. The first two subsections (i.e., speech tokenizer and model architecture) provide a practical base framework for CPT, where the third section discusses the novel CPT strategy applied to the model training.

3.1 Speech Tokenizer

As discussed in Sec. 2, we utilize neural codecs as the speech tokenizer in this work. The speech tokenizer, denoted as $\text{Tokenizer}(\cdot)$, consists of an $\text{Encoder}(\cdot)$, a quantizer $\text{Quantizer}(\cdot)$, and a de-

coder $\text{Decoder}(\cdot)$. The quantizer includes a set of L codebooks, where the i^{th} codebook, $\mathcal{B}^i = \{b_1^i, b_2^i, \dots, b_{B^i}^i\}$, contains B^i codes.

Given a sampled speech signal $S \in \mathbb{R}^{1 \times T_S}$ with length T_S , the encoder processes the signal into hidden states $Q = \text{Encoder}(S) \in \mathbb{R}^{D^{\text{embed}} \times T_Q}$, where D^{embed} represents the dimension of each frame in the hidden states, and T_Q is the temporal sequence length. The quantizer then transforms Q into discrete codes $C \in (\mathcal{B}^1 \times \mathcal{B}^2 \times \dots \times \mathcal{B}^L)^{T_C}$ across T_C frames:

$$C, E = \text{Quantizer}(Q | \mathcal{B}^1, \mathcal{B}^2, \dots, \mathcal{B}^L), \quad (1)$$

where the hidden states $E \in \mathbb{R}^{D^{\text{embed}} \times T_C}$ is constructed from the discrete codes C . Specifically, using the codebooks, the discrete code C is firstly mapped into corresponding embeddings $M \in \mathbb{R}^{L \times D^{\text{embed}} \times T_C}$. M are then compressed along the codebooks’ axis to form E . To reconstruct the speech signal S , the decoder takes E as input and generates the reconstructed signal $\hat{S} = \text{Decoder}(E)$. The discrete codes C are used as the I/O of speech signals for codec-based speech LLM.

3.2 Model Architecture

The model architecture is illustrated in Figure 1. In the first step, both speech codec tokens and textual tokens are transformed into multi-modal shared embeddings. For codec tokens, we select a subset of $L' (\leq L)$ streams of codes, denoted as $C' \in (\mathcal{B}^1 \times \mathcal{B}^2 \times \dots \times \mathcal{B}^{L'})^{T_C}$, from the full set of codec tokens C in the input speech stream.¹ Each code in C' is mapped into corresponding embeddings $M' \in \mathbb{R}^{L' \times D^{\text{embed}} \times T_C}$ and is then aggregated into $E^{\text{speech}} \in \mathbb{R}^{D^{\text{embed}} \times T_C}$. Textual tokens $A \in \mathcal{V}^{T_A}$ with a vocabulary of \mathcal{V} are directly transformed into embeddings $E^{\text{text}} \in \mathbb{R}^{D^{\text{embed}} \times T_A}$ with a sequence length of T_A . Both E^{speech} and E^{text} (denoted as E) are used seamlessly within the LLM’s decoder-only architecture. For simplicity, we denote $R \in (\mathcal{B}^1 \times \mathcal{B}^2 \times \dots \times \mathcal{B}^{L'} \times \mathcal{V})^{T_R}$ that represents both speech codec tokens and text tokens. The T_R can be the length of either a speech codec token sequence (i.e., T_C) or a text token sequence (i.e., T_A).

The core of the LLM is adapted from Qwen1.5, a Transformer-based LLM (Yang et al., 2024a).²

¹ L' is a hyperparameter chosen to balance computational efficiency, modeling complexity in LLM training, and audio reconstruction quality.

²While Qwen1.5 has demonstrated superior performance compared to other recent open-source textual LLMs (Yang et al., 2024a), it also offers a 0.5B version that fits within our

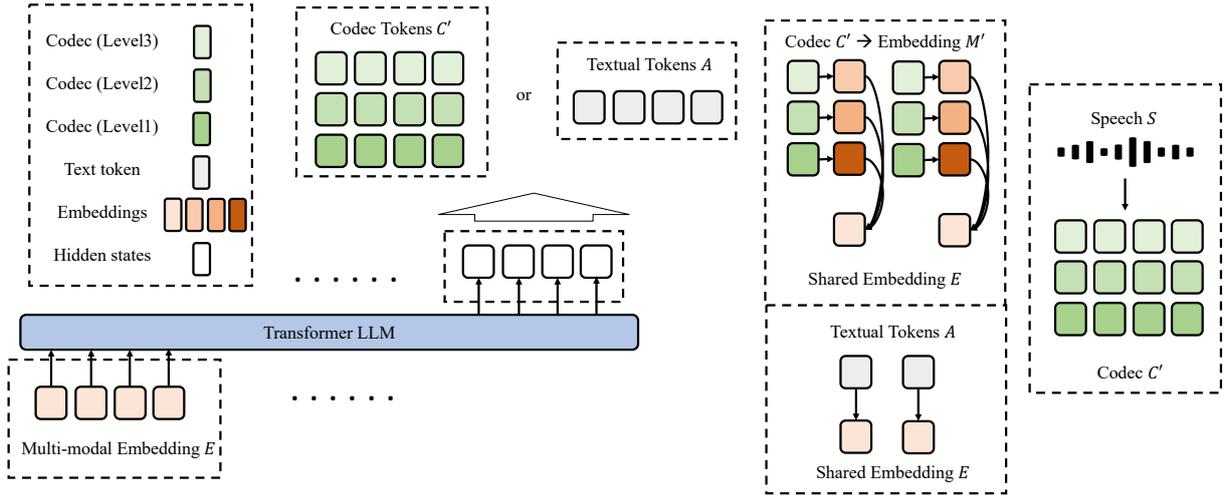


Figure 1: The architecture of the codec-based speech LLM. Codec tokens and textual tokens are converted into multi-modal shared embeddings, as shown on the right side of the figure. These shared embeddings are then fed into a Transformer-based LLM, which features parallel prediction heads designed for predicting either codec tokens or textual tokens. Details are discussed in Sec. 3.2.

To retain the majority of pre-trained knowledge for CPT, we preserve the original Transformer architecture along with the same textual tokenizer and textual token embeddings. Given the hidden states from the final Transformer layer, the language model prediction head is partially initialized with the pre-trained Qwen1.5 weights, while the portion related to the extended vocabulary for codec tokens is initialized randomly. Additional language prediction heads are introduced to predict subsequent codec levels *in parallel*. Unlike recent works (Chen et al., 2024a; Copet et al., 2023; Yang et al., 2024b), we do not enforce dependency constraints between codec predictions. Compared to approaches that use non-autoregressive networks (Wang et al., 2023) or multi-scale Transformers (Yang et al., 2024b), our parallel prediction style remains simple yet efficient during inference.

3.3 Model Training

The formulation discussed in the previous two sections provides a unified interface for speech and text modalities within a single end-to-end model, where speech is represented in multiple streams and text in a single stream. Notably, this is the first approach to seamlessly integrate multi-stream speech codecs with text tokens.

However, this formulation has limitations when incorporating a pre-trained textual LLM due to significant modality mismatches, particularly with speech codecs. To address this, we propose us-

computational constraints. Therefore, in this work, we focus on Qwen1.5 as our base textual LLM.

ing CPT to align modalities, enabling a unified model for both understanding and generation. This approach leverages the strong foundation of textual LLMs while maintaining high-fidelity speech acoustics through advanced codecs.

In the original textual LLM, the pre-training objective is defined as:

$$\theta_{\text{PT-A}} = \operatorname{argmax}_{\theta} \prod_t^{T_A} P_{\theta}(a_t | a_0, a_1, \dots, a_{t-1}), \quad (2)$$

where θ is the language model parameters. When applying CPT using only speech data, the objective transforms into:

$$\theta_{\text{CPT-C}} = \operatorname{argmax}_{\theta} \prod_t^{T_C} P_{\theta}(c_t | c_0, c_1, \dots, c_{t-1}). \quad (3)$$

The purpose of CPT is to inject knowledge about speech modality into the LLM transformer base. Therefore, we conduct two pre-training configurations: one with speech data only following Eq. (3), and the other with bi-modal (i.e., speech and text) datasets. Details of data preparation are discussed in Sec. 4.1. We denote a sequence of speech and text tokens $R = (r_1, r_2, \dots, r_{T_R})$ with a length of T_R . $r_t \in (\mathcal{B}^1 \times \mathcal{B}^2 \times \dots \times \mathcal{B}^{L'} \times \mathcal{V})$ can represent either a speech codec token $(c_t^1, c_t^2, \dots, c_t^{L'}, \text{null})$ or text token $(\text{null}, \dots, \text{null}, a_t)$. The pre-training objective follows the basic next-token prediction task:

$$\theta_{\text{CPT}} = \operatorname{argmax}_{\theta} \prod_t^{T_R} P_{\theta}(r_t | r_0, r_1, \dots, r_{t-1}), \quad (4)$$

where θ is the model illustrated in Fig. 1 and r_0 is the start of sentence.

Following the CPT stage, we conduct SFT where a prompt token sequence R^{inp} is first processed before the target sequence R^{tgt} . The objective for the SFT is:

$$\theta_{\text{SFT}} = \operatorname{argmax}_{\theta} P_{\theta}(R^{\text{tgt}} | R^{\text{inp}}). \quad (5)$$

4 Experimental Setup

In this section, we detail the experimental settings for CPT and SFT to evaluate the effectiveness of the proposed method across various downstream tasks. Specifically, we compare CPT-ed models with those that are either randomly initialized or initialized from a textual LLM.

The basic speech tokenizer follows the Sound-Stream architecture (Zeghidour et al., 2021). Instead of using a complex short-term Fourier transform (STFT) discriminator, we adopt the discriminators from (Kumar et al., 2023), which include a multi-frequency STFT discriminator, a multi-scale discriminator, and a multi-period discriminator. For more detailed hyperparameters of the speech tokenizer, please refer to Appendix A.

4.1 Pre-training Data Preparation

As mentioned in Sec. 1, we focus on codec-based speech LLMs. Therefore, in the following experiments, we primarily consider bilingual scenarios, using both English and Mandarin data.

The pre-training data consists of around 140k hours of English and Mandarin speech, along with corresponding transcriptions and translations. Given the limited availability of open-source Mandarin speech data, we incorporate 70k hours of in-house Mandarin data to maintain a balance between English and Mandarin in the pre-training corpus. More detailed information about the corpora used for pre-training is provided in Appendix B.1. For training consistency on punctuation, we use a BERT-based punctuation restoration model³ to recover the punctuation for textual data without punctuation (e.g., Librispeech (Panayotov et al., 2015)).

Most of the corpora included were originally designed for ASR or TTS purposes and did not include translations in their official release, which may prevent full alignment of the semantic spaces between the two languages. To semantically align

³<https://huggingface.co/felflare/bert-restore-punctuation>

English and Mandarin in the speech LLM, we supplement the data using an internal machine translation model to generate translations between English and Mandarin (EN->ZH and ZH->EN) based on the original speech transcriptions.⁴ For each utterance in the pre-training dataset, we use the paired tuple (i.e., speech, transcription, and translation) to create six tasks: (1) speech continuation, which predicts future speech tokens based on previous segments; (2) language modeling, which generates the following textual tokens from a given context; (3) ASR, which transcribes speech into text; (4) TTS, which generates speech from text using a target speaker’s speech segment as a prompt; (5) speech-to-text translation (S2TT), which translates speech in one language to text in another; and (6) text-to-speech translation (T2ST), which generates translated speech from a source text using a target speaker’s prompt.

We did not include S2ST during continual pre-training because (1) our ASR/TTS corpora lack native S2ST pairs, making high-quality paired data scarce, and (2) synthesizing speech that fully preserves speaker identity and natural prosody is challenging. Instead, we reserve S2ST as an unseen task and later verify its performance via supervised fine-tuning, allowing us to assess whether our CPT framework has successfully aligned the English and Mandarin semantic spaces.

Based on the six tasks mentioned, the final sequence is formulated using the pattern “<Condition><Prompt><Target>”. Additionally, a boundary token is inserted for each speech or text segment to mark the start or end of the corresponding segment. Depending on the specific task type, we enforce slightly different policies generating the data. The specifics are detailed in Appendix B.2. Noted that the natural language prompts are only in textual modality.

4.2 Continual Pre-training

We adopt Megatron-LM⁵ with tensor parallelism as our training framework (Shoeybi et al., 2019), using Qwen1.5-0.5B as the base pre-trained textual LLM (Yang et al., 2024a).

As formulated in Eq. (4), we carry out two CPT experiments: one using only speech data, and the other incorporating both speech and text modalities.

⁴We use our internal system because it’s fine-tuned on our speech corpora and integrated with our data pipeline, ensuring more consistent and accurate translations for our tasks.

⁵<https://github.com/NVIDIA/Megatron-LM>

425 For the mixed-modality training, we aim to balance
426 different domains and tasks by constructing a data
427 loader that samples data from various tasks. Specif-
428 ically, the loader randomly selects data for each
429 of the six speech-related tasks with a 15% proba-
430 bility. To maintain the reasoning capacity of the
431 original textual LLM and avoid catastrophic forget-
432 ting, we include text-only data for the remaining
433 10%. Of this textual data, 5% comes from gen-
434 eral domains such as books, YouTube titles, and
435 Wikipedia, while the other 5% is sourced from an
436 in-house machine translation (MT) corpus to en-
437 hance the model’s translation capabilities. This de-
438 sign aligns with previous CPT approaches in both
439 textual and multimodal LLMs, where even limited
440 text-based pre-training acts as a stabilizing regular-
441 izer, anchoring the model’s internal representations
442 to those learned during pre-training (Zhai et al.,
443 2023; Sun et al., 2020). Detailed training hyperpa-
444 rameters are discussed in Appendix C.

445 Additionally, to identify the effect of CPT, we
446 conduct experiments with a random-initialized
447 0.5B model that has the same architecture as
448 Qwen1.5-0.5B and the pre-trained Qwen1.5-0.5B
449 for comparison.⁶

4.3 Downstream Tasks

451 As formulated in Eq. (5), the downstream evalua-
452 tion is performed by fine-tuning the pre-trained
453 model across four tasks: ASR, TTS, S2TT, and
454 S2ST. ASR, TTS, and S2TT are used to assess the
455 model’s basic speech understanding and generation
456 abilities. S2ST can be broken down into ASR, MT,
457 and TTS, making it a more comprehensive task
458 that evaluates both understanding and generation
459 capabilities.

460 All four fine-tuning tasks follow the sequence
461 formulation discussed in Sec. 4.1, with different
462 *<Condition>*, *<Prompt>*, and *<Target>*, related
463 to the tasks. Notably, the task prompts are also
464 generated following the pipeline in Sec. 4.1 but
465 they do not overlap with the pre-training prompts.
466 **ASR:** For the ASR task, we evaluate the model
467 using the Librispeech dataset for English and the
468 Aishell2 dataset for Mandarin (Panayotov et al.,
469 2015; Du et al., 2018). The *<Condition>* consists
470 of a sequence of speech codec tokens, followed
471 by a natural language prompt. As described in
472 Sec. 4.1, the *<Target>* sequence is the transcrip-
473 tion with restored punctuation. To further improve

⁶They are denoted as “No Initialization” and “Text LLM Initialization” in the following discussion.

474 model performance, we apply random time-domain
475 masking, similar to the approach in (Chang et al.,
476 2023, 2024c).

477 During inference, the model autoregressively
478 predicts the transcription by feeding the *<Con-
479 dition><Prompt>* sequence into the system. To
480 enhance decoding performance, we employ beam
481 search, using a beam size of 8. We measure word
482 error rate (WER) for English ASR and character
483 error rate (CER) for Mandarin ASR.

484 **TTS:** We focus on the multi-speaker TTS task with
485 LibriTTS (Zen et al., 2019). For the task formu-
486 lation, we follow VaLL-E style input where the
487 condition includes a three-second speaker prompt
488 and the text (Wang et al., 2023). For the target text,
489 we restore the punctuation similar to the ASR task.

490 Since greedy search tends to produce trivial out-
491 puts, we adopt a sampling-based inference with
492 a top- k strategy, setting $k = 30$, as used in prior
493 works (Wang et al., 2023; Yang et al., 2024b; Tian
494 et al., 2024). To further increase the diversity of
495 generated speech, we re-scale the predicted logits
496 using a temperature of 1.5.

497 For evaluation, we use: WER from a pre-trained
498 ASR model, speaker similarity (SPK-SIM) from
499 a pre-trained speaker embedding model, and an
500 automatic speech quality predictor based on a
501 pre-trained mean opinion score (MOS) predictor.
502 Specifically, we use Whisper-large-V3 (Radford
503 et al., 2023) for WER evaluation, a pre-trained
504 Rawnet3 model (Jung et al., 2024) trained on Vox-
505 celeb for speaker embedding extraction, and UT-
506 MOS (Saeki et al., 2022) as the MOS predictor. Fol-
507 lowing common practice in previous works (Wang
508 et al., 2023; Yang et al., 2024b; He et al., 2024;
509 Tian et al., 2024), we generate five samples per test
510 instance using the sampling strategy and report the
511 average score across each metric.

512 **S2TT:** The S2TT adopts a task formulation as the
513 ASR task by simply replacing the input speech with
514 source language speech and target transcription
515 with target language translation. The prompts are
516 changed to task-related prompts accordingly. We
517 test two corpora: CoVOST2 and GigaST (Wang
518 et al., 2021; Ye et al., 2023). For CoVOST2,
519 we focus on two translation directions, including
520 English-to-Mandarin (EN->ZH) and Mandarin-to-
521 English (ZH->EN). For GigaST, we only focus on
522 EN->ZH. We use SacreBLEU to evaluate the pre-
523 diction results with BLEU score (Post, 2018).

524 **S2ST:** We conduct S2ST using the GigaS2S cor-
525 pus, which supplements the GigaST corpus with a

Table 1: ASR performance on LibriSpeech and Aishell2. Models marked with * indicate pre-trained models did not undergo continual pre-training. + stands models that do not use neural codecs as their speech representation. We report WER for Librispeech and CER for Aishell2.

Models	Param.	LibriSpeech		Aishell2
		Test-clean	Test-other	Test-overall
VoxLM* (Maiti et al., 2024)	1B	2.7	6.5	-
AnyGPT (Zhan et al., 2024)	7B	8.5	-	-
No Initialization*		5.5	9.5	15.5
Text LLM Initialization*		4.8	8.5	13.1
Speech CPT	1B	5.5	8.9	13.0
Speech & Text CPT		3.7	6.3	7.2

single-speaker TTS model (Ye et al., 2023).⁷ Due to data constraints, we focus only on the English-to-Mandarin (EN->ZH) translation direction. Since the target speech is single-speaker, the *<Condition>* consists solely of source language speech (i.e., English). The prompt generation process follows the same approach as in previous tasks. Consistent with the method used in AudioPalm (Rubenstein et al., 2023), we directly assign the *<Target>* as the codec sequence of the target speech.

For inference, we use the same top-*k* sampling strategy employed in the TTS task. For evaluation, we measure translation quality using ASR-BLEU and speech quality using UTMOS.⁸

5 Results and Discussion

5.1 ASR

The experimental results for ASR are presented in Table 1. The best performance is achieved by the model that underwent CPT with both speech and text modalities. This outcome is expected since ASR is one of the tasks included in the pre-training phase. The model that received only speech CPT performed worse on Librispeech, indicating that focusing solely on speech continuation within the speech codec pre-training does not necessarily enhance speech recognition performance.

We also compare these results with reference performances from other speech LLM-based models on Librispeech. VoxLM, which uses SSL-based tokens as its modeling unit (Maiti et al., 2024), and AnyGPT, a multi-modal LLM that uses SpeechTokenizer—a speech neural codec that additionally

⁷<https://github.com/SpeechTranslation/GigaS2S>

⁸Although UTMOS was trained on English speech, which might introduce some language mismatch in the scoring, prior work (Huang et al., 2022) has shown that UTMOS still achieved reasonable correlation scores when evaluating out-of-domain Chinese speech. Thus, we continue to use it for our speech quality evaluation.

Table 2: TTS performance on LibriTTS. Models marked with * indicate pre-trained models that did not undergo continual pre-training. ° corresponds to a version trained on LibriTTS.

Models	UTMOS	WER	SPK-SIM
UniAudio° (Yang et al., 2024b)	3.64	13.1	0.43
No Initialization*	3.01	17.5	0.55
Text LLM Initialization*	2.78	18.8	0.51
Speech CPT	3.65	3.7	0.66
Speech & Text CPT	3.59	3.7	0.65

distills representations from speech self-supervised models (Zhan et al., 2024; Zhang et al., 2024)—are included in the comparison. While the proposed model with CPT slightly degraded from VoxLM on the Librispeech text-clean set, it outperforms on the more realistic test-other set.⁹ This potentially suggests that CPT successfully improves the model’s understanding ability in codec-based speech LLMs, bringing its performance closer to that of SSL representations, which are known to excel in understanding tasks.

5.2 TTS

The TTS performance results are presented in Table 2. The best-performing system is the model that uses speech-only CPT, indicating that speech continuation pre-training can significantly enhance speech generation quality and ease the challenges associated with speech generation modeling. The model with joint speech and text CPT achieves comparable performance in terms of intelligibility, as measured by WER, and in speaker prompt understanding, as indicated by SPK-SIM. Overall, models that apply CPT have demonstrated superior performance compared to those with or without textual LLM initialization.

As shown in Table 2, We also conduct experiments on the LibriTTS dataset using the UniAudio-based model (i.e., multi-scale transformer-based language model TTS model) in ESPnet (Shi et al., 2024; Tian et al., 2024; Yang et al., 2024b). The same speech tokenizer as our CPT-ed model is employed for these experiments. Compared to the TTS-specialized model, the results demonstrate that the proposed method generates speech with comparable quality, as measured by UTMOS, while significantly improving intelligibility and speaker style transfer, as reflected by much lower WER and higher SPK-SIM scores.

⁹Due to differences in pre-training data and model size, concrete comparisons are challenging.

Table 3: S2TT performance on CoVOST2 and GigaST. The performance is reported in BLEU. Models marked with * indicate pre-trained models that did not undergo continual pre-training. • stands that external machine translation data is used.

Pre-training	CoVoST2		GigaST
	EN -> ZH	ZH -> EN	EN -> ZH
Fairseq ST (Wang et al., 2021)	25.4	5.8	-
OWSM-v3 (Peng et al., 2023)	33.4	13.6	-
GigaST* (Ye et al., 2023)	-	-	38.0
LLM-ST* (Huang et al., 2023)	-	-	39.6
No Initialization*	25.5	5.8	30.4
Text LLM Initialization*	28.9	9.9	33.2
Speech CPT	24.8	5.4	33.1
Speech & Text CPT	33.1	16.1	37.5

Combining these results with those from Sec. 5.1, our proposed method not only enhances the ASR performance of codec-based speech LLMs but also maintains high TTS quality, ensuring no degradation in speech generation compared to other specialized codec-based speech language models without textual LLM initialization.

5.3 S2TT

The results for speech-to-text translation are shown in Table 3. In both the CoVOST2 and GigaST datasets, models with CPT using both speech and text modalities demonstrate significant improvements in the S2TT task, highlighting their effectiveness in understanding tasks. Notably, the model with speech-only CPT performs worse than models without initialization. This result aligns with the ASR findings but contrasts with the TTS results, indicating that a focus on speech generation does not necessarily enhance speech understanding ability.

We also present results from related works (Wang et al., 2021; Peng et al., 2023; Ye et al., 2023; Huang et al., 2023). Our proposed model shows better performance than Fairseq-ST and achieves performance comparable to OWSM-v3, which focuses on understanding tasks. For GigaST, since both models incorporate external machine translation data, a direct comparison is not possible. But we observe that the model with joint-modality CPT has achieved performance similar to these models that are specifically designed for S2TT tasks.

5.4 S2ST

The results for S2ST are presented in Table 4. The best performance is achieved by the model with joint speech-text CPT. Notably, unlike ASR, TTS, and S2TT, S2ST is not included as a task in the CPT phase. However, both models that underwent

Table 4: S2ST performance on GigaST. † indicates that the ASR-BLEU scores were calculated using different ASR systems, as described in Sec. 5.4.

Pre-training	GigaST	
	ASR-BLEU	UTMOS
Vec-Tok† (Zhu et al., 2023)	21.6	-
HW-TSC† (Wu et al., 2024d)	33.6	-
Speech CPT	28.0	3.41
Speech & Text CPT	33.4	3.66

CPT still demonstrated strong S2ST performance. In contrast, models without CPT, even after extensive hyper-parameter tuning, struggled to converge effectively on the S2ST task.¹⁰

For comparison, we also include results from two prior works (Zhu et al., 2023; Wu et al., 2024d). Vec-Tok uses an end-to-end architecture with additional emphasis on source speaker style transfer, while the HW-TSC S2ST model is built using a cascaded ASR and MT system. It is important to note that the results are not directly comparable due to differences in the ASR models used for evaluation. However, we observe that the models with CPT show performance potentially comparable to the cascaded approach, highlighting the effectiveness of CPT for this task.

6 Conclusion

We explore continual pre-training as an effective strategy to extend codec-based speech LLMs for speech translation-related tasks. By carefully formulating our pre-training data, we adapt a pre-trained textual LLM in two configurations—one with speech-only data and another with a joint speech-text approach. Our extensive experiments on ASR, TTS, S2TT, and S2ST tasks show that continual pre-training can significantly enhance performance. In particular, speech-only continual pre-training yields notable improvements for TTS, while joint speech-text continual pre-training strikes a balance between understanding and generation, ultimately delivering high-quality end-to-end S2ST. These findings underscore the potential of continual pre-training in addressing issues such as catastrophic forgetting and modality mismatch, thereby advancing the development of robust multimodal language models.¹¹

¹⁰Due to the non-convergence, we did not put the results in Table 4.

¹¹Some generated audio samples are available at <https://hiddenmeprivate.github.io/>

7 Limitations

We demonstrate through several experiments that CPT can potentially help balance the generation and understanding abilities in codec-based speech LLMs. However, we acknowledge the following limitations:

- **Limited model size:** Due to computational and time constraints, we were unable to train a larger-scale model. In comparison to many recent models with 7B or more parameters, the smaller capacity of our model may limit the full exploration of our proposed approach.
- **Difficulty in comparison to recent works:** Performing fair comparisons with recent models is challenging for several reasons, such as mismatched pre-trained textual LLMs (which we could not utilize due to computational limitations), mismatched datasets, and the difficulty of reproducing results with limited resources. As a result, most reference models cannot be directly compared to our models.
- **Use of in-house data:** For continual pre-training, we employed some in-house data to reduce bias present in open-source data (e.g., language distribution, task distribution, and demographic distribution). However, due to privacy agreements and licensing issues, this data cannot be openly shared.
- **Comprehensiveness of ablation studies:** Given our computational budgets, we were unable to conduct full-cycle ablation studies for all aspects of the proposed methodology. Instead, we modified and tested design choices during early-stage training, where performance may not be fully representative of the later-stage training results.

8 Ethical Statement and Potential Risks

The development of codec-based speech language models through CPT has the potential to enhance a wide range of speech-related tasks, including speech-to-text and text-to-speech translations, thereby contributing to advancements in communication technologies. However, the use of these models must be approached with caution due to several ethical considerations:

Speech data often contains sensitive personal information. The use of speech language models necessitates strict adherence to privacy laws

and regulations, such as GDPR, to ensure that personally identifiable information is not inadvertently exposed or misused. Robust mechanisms for data anonymization and secure storage must be employed to prevent unauthorized access or exploitation of individuals' speech data.

The ability to generate high-quality speech through TTS or S2ST poses risks related to the generation of deepfake audio or other forms of speech-based manipulation. Misuse of this technology could potentially lead to the spread of misinformation or impersonation, raising concerns over its potential role in fraudulent activities or the dissemination of false information. Measures such as watermarking or other forms of verifiable speech generation could be explored to mitigate these risks.

In developing codec-based speech language models, we commit to adhering to ethical principles that prioritize user privacy, data security, and fairness, while also actively working to mitigate any negative societal impacts that may arise from misuse of this technology.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Philip Anastassiou, Jiawei Chen, Jitong Chen, Yuanzhe Chen, Zhuo Chen, Ziyi Chen, Jian Cong, Lelai Deng, Chuang Ding, Lu Gao, et al. 2024. Seed-TTS: A family of high-quality versatile speech generation models. *arXiv preprint arXiv:2406.02430*.
- Rosana Ardila, Megan Branson, Kelly Davis, Michael Kohler, Josh Meyer, Michael Henretty, Reuben Morais, Lindsay Saunders, Francis Tyers, and Gregor Weber. 2020. **Common voice: A massively-multilingual speech corpus**. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4218–4222, Marseille, France. European Language Resources Association.
- Hui Bu, Jiayu Du, Xingyu Na, Bengu Wu, and Hao Zheng. 2017. Aishell-1: An open-source mandarin speech corpus and a speech recognition baseline. In *2017 20th conference of the oriental chapter of the international coordinating committee on speech databases and speech I/O systems and assessment (O-COCOSDA)*, pages 1–5. IEEE.
- Kai-Wei Chang, Haibin Wu, Yu-Kai Wang, Yuan-Kuei Wu, Hua Shen, Wei-Cheng Tseng, Iu-thing Kang, Shang-Wen Li, and Hung-yi Lee. 2024a. Speech-prompt: Prompting speech language models for

768	speech processing tasks. <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> .	824
769		825
770	Xuankai Chang, Jiatong Shi, Jinchuan Tian, Yuning Wu, Yuxun Tang, Yihan Wu, Shinji Watanabe, Yossi Adi, Xie Chen, and Qin Jin. 2024b. The interspeech 2024 challenge on speech processing using discrete units . In <i>Interspeech 2024</i> , pages 2559–2563.	826
771		827
772		828
773		829
774		830
775	Xuankai Chang, Brian Yan, Kwanghee Choi, Jee-Weon Jung, Yichen Lu, Soumi Maiti, Roshan Sharma, Jiatong Shi, Jinchuan Tian, Shinji Watanabe, et al. 2024c. Exploring speech recognition, translation, and understanding with discrete speech units: A comparative study. In <i>ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 11481–11485. IEEE.	831
776		832
777		833
778		834
779		835
780		836
781		837
782		838
783	Xuankai Chang, Brian Yan, Yuya Fujita, Takashi Maekaku, and Shinji Watanabe. 2023. Exploration of efficient end-to-end ASR using discretized input from self-supervised learning . In <i>INTERSPEECH 2023</i> , pages 1399–1403.	839
784		840
785		841
786		842
787		843
788	Guoguo Chen, Shuzhou Chai, Guan-Bo Wang, Jiayu Du, Wei-Qiang Zhang, Chao Weng, Dan Su, Daniel Povey, Jan Trmal, Junbo Zhang, Mingjie Jin, Sanjeev Khudanpur, Shinji Watanabe, Shuaijiang Zhao, Wei Zou, Xiangang Li, Xuchen Yao, Yongqing Wang, Zhao You, and Zhiyong Yan. 2021. Gigaspeech: An evolving, multi-domain asr corpus with 10,000 hours of transcribed audio . In <i>Interspeech 2021</i> , pages 3670–3674.	844
789		845
790		846
791		847
792		848
793		849
794		850
795		851
796		852
797	Sanyuan Chen, Shujie Liu, Long Zhou, Yanqing Liu, Xu Tan, Jinyu Li, Sheng Zhao, Yao Qian, and Furu Wei. 2024a. VALL-E 2: Neural codec language models are human parity zero-shot text to speech synthesizers. <i>arXiv preprint arXiv:2406.05370</i> .	853
798		854
799		855
800		856
801		857
802	Zhehuai Chen, He Huang, Andrei Andrusenko, Oleksii Hrinchuk, Krishna C Puvvada, Jason Li, Subhankar Ghosh, Jagadeesh Balam, and Boris Ginsburg. 2024b. SALM: Speech-augmented language model with in-context learning for speech recognition and translation. In <i>ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 13521–13525. IEEE.	858
803		859
804		860
805		861
806		862
807		863
808		864
809		865
810	Zhehuai Chen, Yu Zhang, Andrew Rosenberg, Bhuvana Ramabhadran, Pedro J. Moreno, Ankur Bapna, and Heiga Zen. 2022. MAESTRO: Matched speech text representations through modality matching . In <i>Interspeech 2022</i> , pages 4093–4097.	866
811		867
812		868
813		869
814		870
815	Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, Chang Zhou, and Jingren Zhou. 2024a. Qwen2-audio technical report . <i>Preprint</i> , arXiv:2407.10759.	871
816		872
817		873
818		874
819		875
820	Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, et al. 2024b. Qwen2-audio technical report . <i>arXiv preprint arXiv:2407.10759</i> .	876
821		877
822		878
823		879
	Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. 2023. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models . <i>Preprint</i> , arXiv:2311.07919.	880
		881
	Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre Defossez. 2023. Simple and controllable music generation . In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 47704–47720. Curran Associates, Inc.	882
		883
	Alexandre Défossez, Laurent Mazaré, Manu Orsini, Amélie Royer, Patrick Pérez, Hervé Jégou, Edouard Grave, and Neil Zeghidour. 2024a. Moshi: a speech-text foundation model for real-time dialogue .	884
		885
	Alexandre Défossez, Laurent Mazaré, Manu Orsini, Amélie Royer, Patrick Pérez, Hervé Jégou, Edouard Grave, and Neil Zeghidour. 2024b. Moshi: a speech-text foundation model for real-time dialogue . <i>arXiv preprint arXiv:2410.00037</i> .	886
		887
	Kunal Dhawan, Nithin Rao Koluguri, Ante Jukić, Ryan Langman, Jagadeesh Balam, and Boris Ginsburg. 2024. Codec-asr: Training performant automatic speech recognition systems with discrete speech representations . In <i>Interspeech 2024</i> , pages 2574–2578.	888
		889
	Jiayu Du, Xingyu Na, Xuechen Liu, and Hui Bu. 2018. Aishell-2: Transforming mandarin asr research into industrial scale . <i>arXiv preprint arXiv:1808.10583</i> .	890
		891
	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The Llama 3 herd of models . <i>arXiv preprint arXiv:2407.21783</i> .	892
		893
	Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text for language modeling . <i>arXiv preprint arXiv:2101.00027</i> .	894
		895
	Yuan Gong, Alexander H Liu, Hongyin Luo, Leonid Karlinsky, and James Glass. 2023. Joint audio and speech understanding . In <i>2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)</i> , pages 1–8. IEEE.	896
		897
	Yuan Gong, Hongyin Luo, Alexander H. Liu, Leonid Karlinsky, and James R. Glass. 2024. Listen, think, and understand . In <i>The Twelfth International Conference on Learning Representations</i> .	898
		899
	Yiwei Guo, Zhihan Li, Hankun Wang, Bohan Li, Chongtian Shao, Hanglei Zhang, Chenpeng Du, Xie Chen, Shujie Liu, and Kai Yu. 2025. Recent advances in discrete speech tokens: A review . <i>arXiv preprint arXiv:2502.06490</i> .	900
		901
	Chi Han, Mingxuan Wang, Heng Ji, and Lei Li. 2021. Learning shared semantic space for speech-to-text	902

878		translation . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 2214–2225, Online. Association for Computational Linguistics.	
879			
880			
881			
882	Haorui He, Zengqiang Shang, Chaoren Wang, Xuyuan Li, Yicheng Gu, Hua Hua, Liwei Liu, Chen Yang, Jiaqi Li, Peiyang Shi, et al. 2024. Emilia: An extensive, multilingual, and diverse speech dataset for large-scale speech generation. <i>arXiv preprint arXiv:2407.05361</i> .		
883			
884			
885			
886			
887			
888	François Hernandez, Vincent Nguyen, Sahar Ghanay, Natalia Tomashenko, and Yannick Esteve. 2018. TED-LIUM 3: Twice as much data and corpus repartition for experiments on speaker adaptation. In <i>Speech and Computer: 20th International Conference, SPECOM 2018</i> , pages 198–208. Springer.		
889			
890			
891			
892			
893			
894	Shujie Hu, Long Zhou, Shujie Liu, Sanyuan Chen, Hongkun Hao, Jing Pan, Xunying Liu, Jinyu Li, Sunit Sivasankaran, Linqun Liu, et al. 2024. Wavllm: Towards robust and adaptive speech large language model. <i>arXiv preprint arXiv:2404.00656</i> .		
895			
896			
897			
898			
899	Chien-yu Huang, Wei-Chih Chen, Shu-wen Yang, Andy T Liu, Chen-An Li, Yu-Xiang Lin, Wei-Cheng Tseng, Anuj Diwan, Yi-Jen Shih, Jiatong Shi, et al. 2024a. Dynamic-SUPERB phase-2: A collaboratively expanding benchmark for measuring the capabilities of spoken language models with 180 tasks. <i>arXiv preprint arXiv:2411.05361</i> .		
900			
901			
902			
903			
904			
905			
906	Chien-yu Huang, Ke-Han Lu, Shih-Heng Wang, Chi-Yuan Hsiao, Chun-Yi Kuan, Haibin Wu, Siddhant Arora, Kai-Wei Chang, Jiatong Shi, Yifan Peng, et al. 2024b. Dynamic-SUPERB: Towards a dynamic, collaborative, and comprehensive instruction-tuning benchmark for speech. In <i>ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 12136–12140. IEEE.		
907			
908			
909			
910			
911			
912			
913			
914			
915	Rongjie Huang, Chunlei Zhang, Yongqi Wang, Dongchao Yang, Jinchuan Tian, Zhenhui Ye, Luping Liu, Zehan Wang, Ziyue Jiang, Xuankai Chang, et al. 2024c. Make-a-voice: Revisiting voice large language models as scalable multilingual and multitask learners. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10929–10942.		
916			
917			
918			
919			
920			
921			
922			
923	Wen Chin Huang, Erica Cooper, Yu Tsao, Hsin-Min Wang, Tomoki Toda, and Junichi Yamagishi. 2022. The VoiceMOS challenge 2022 . In <i>Interspeech 2022</i> , pages 4536–4540.		
924			
925			
926			
927	Zhichao Huang, Rong Ye, Tom Ko, Qianqian Dong, Shanbo Cheng, Mingxuan Wang, and Hang Li. 2023. Speech translation with large language models: An industrial practice. <i>arXiv preprint arXiv:2312.13585</i> .		
928			
929			
930			
931	Jee-Weon Jung, Wangyou Zhang, Jiatong Shi, Zakaria Aldeneh, Takuya Higuchi, Alex Gichamba, Barry John Theobald, Ahmed Hussien Abdelaziz, and Shinji Watanabe. 2024. ESPnet-SPK: full pipeline speaker		
932			
933			
934			
		embedding toolkit with reproducible recipes, self-supervised front-ends, and off-the-shelf models . In <i>Interspeech 2024</i> , pages 4278–4282.	935
			936
			937
	Eugene Kharitonov, Ann Lee, Adam Polyak, Yossi Adi, Jade Copet, Kushal Lakhotia, Tu Anh Nguyen, Morgane Riviere, Abdelrahman Mohamed, Emmanuel Dupoux, et al. 2022. Text-free prosody-aware generative spoken language modeling. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8666–8681.		938
			939
			940
			941
			942
			943
			944
			945
	Jaehyeon Kim, Keon Lee, Seungjun Chung, and Jaewoong Cho. 2024. CLam-TTS: Improving neural codec language model for zero-shot text-to-speech . In <i>The Twelfth International Conference on Learning Representations</i> .		946
			947
			948
			949
			950
	Minje Kim and Jan Skoglund. 2024. Neural speech and audio coding . <i>Preprint</i> , arXiv:2408.06954.		951
			952
	Vijay Anand Korthikanti, Jared Casper, Sangkug Lym, Lawrence McAfee, Michael Andersch, Mohammad Shoeybi, and Bryan Catanzaro. 2023. Reducing activation recomputation in large transformer models. <i>Proceedings of Machine Learning and Systems</i> , 5:341–353.		953
			954
			955
			956
			957
			958
	Chun-Yi Kuan, Chih-Kai Yang, Wei-Ping Huang, Ke-Han Lu, and Hung-yi Lee. 2024. Speech-Copilot: Leveraging large language models for speech processing via task decomposition, modularization, and program generation . <i>arXiv preprint arXiv:2407.09886</i> .		959
			960
			961
			962
			963
	Rithesh Kumar, Prem Seetharaman, Alejandro Luebs, Ishaan Kumar, and Kundan Kumar. 2023. High-fidelity audio compression with improved rvqgan . <i>ArXiv</i> , abs/2306.06546.		964
			965
			966
			967
	Ann Lee, Peng-Jen Chen, Changhan Wang, Jiatao Gu, Sravya Popuri, Xutai Ma, Adam Polyak, Yossi Adi, Qing He, Yun Tang, Juan Pino, and Wei-Ning Hsu. 2022. Direct speech-to-speech translation with discrete units . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3327–3339, Dublin, Ireland. Association for Computational Linguistics.		968
			969
			970
			971
			972
			973
			974
			975
	Soumi Maiti, Yifan Peng, Shukjae Choi, Jee-weon Jung, Xuankai Chang, and Shinji Watanabe. 2024. VoxLM: Unified decoder-only models for consolidating speech recognition, synthesis and speech, text continuation tasks . In <i>ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 13326–13330. IEEE.		976
			977
			978
			979
			980
			981
			982
	Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. Librispeech: an asr corpus based on public domain audio books. In <i>2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)</i> , pages 5206–5210. IEEE.		983
			984
			985
			986
			987
			988
	Yifan Peng, Jinchuan Tian, Brian Yan, Dan Berrebbi, Xuankai Chang, Xinjian Li, Jiatong Shi, Siddhant		989
			990

991	Arora, William Chen, Roshan Sharma, et al. 2023. Reproducing Whisper-style training using an open-source toolkit and publicly available data. In <i>2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)</i> , pages 1–8. IEEE.	1046
992		1047
993		1048
994		1049
995		1050
996	Matt Post. 2018. A call for clarity in reporting BLEU scores . In <i>Proceedings of the Third Conference on Machine Translation: Research Papers</i> , pages 186–191, Brussels, Belgium. Association for Computational Linguistics.	1051
997		1052
998		1053
999		1054
1000		1055
1001	Vineel Pratap, Qiantong Xu, Anuroop Sriram, Gabriel Synnaeve, and Ronan Collobert. 2020. MLS: A large-scale multilingual dataset for speech research . In <i>Interspeech 2020</i> , pages 2757–2761.	1056
1002		1057
1003		1058
1004		1059
1005	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In <i>International conference on machine learning</i> , pages 28492–28518. PMLR.	1060
1006		1061
1007		1062
1008		1063
1009		1064
1010	J Rottland, Ch Neukirchen, and D Willett. 1997. The WSJ speech database.	1065
1011		1066
1012	Paul K Rubenstein, Chulayuth Asawaroengchai, Duc Dung Nguyen, Ankur Bapna, Zalán Borsos, Félix de Chaumont Quitry, Peter Chen, Dalia El Badawy, Wei Han, Eugene Kharitonov, et al. 2023. AudioPaLM: A large language model that can speak and listen. <i>arXiv preprint arXiv:2306.12925</i> .	1067
1013		1068
1014		1069
1015		1070
1016		1071
1017		1072
1018	Takaaki Saeki, Detai Xin, Wataru Nakata, Tomoki Koriyama, Shinnosuke Takamichi, and Hiroshi Saruwatari. 2022. UTMOS: Utokyo-sarulab system for voicemos challenge 2022 . In <i>Interspeech 2022</i> , pages 4521–4525.	1073
1019		1074
1020		1075
1021		1076
1022		1077
1023	Jiatong Shi, Yun Tang, Ann Lee, Hirofumi Inaguma, Changhan Wang, Juan Pino, and Shinji Watanabe. 2023. Enhancing speech-to-speech translation with multiple TTS targets . In <i>ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 1–5.	1078
1024		1079
1025		1080
1026		1081
1027		1082
1028		1083
1029	Jiatong Shi, Jinchuan Tian, Yihan Wu, Jee weon Jung, Jia Qi Yip, Yoshiki Masuyama, William Chen, Yuning Wu, Yuxun Tang, Massa Baali, Dareen Alharhi, Dong Zhang, Ruifan Deng, Tejes Srivastava, Haibin Wu, Alexander H. Liu, Bhiksha Raj, Qin Jin, Ruihua Song, and Shinji Watanabe. 2024. ESPnet-Codec: Comprehensive training and evaluation of neural codecs for audio, music, and speech . <i>Preprint</i> , arXiv:2409.15897.	1084
1030		1085
1031		1086
1032		1087
1033		1088
1034		1089
1035		1090
1036		1091
1037		1092
1038	Yao Shi, Hui Bu, Xin Xu, Shaoji Zhang, and Ming Li. 2021. Aishell-3: A multi-speaker mandarin tts corpus . In <i>Interspeech 2021</i> , pages 2756–2760.	1093
1039		1094
1040		1095
1041	Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. 2019. Megatron-lm: Training multi-billion parameter language models using model parallelism. <i>arXiv preprint arXiv:1909.08053</i> .	1096
1042		1097
1043		1098
1044		1099
1045		1099
	Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Hao Tian, Hua Wu, and Haifeng Wang. 2020. Ernie 2.0: A continual pre-training framework for language understanding. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 34, pages 8968–8975.	1046
		1047
		1048
		1049
		1050
	Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun MA, and Chao Zhang. 2024a. SALMONN: Towards generic hearing abilities for large language models . In <i>The Twelfth International Conference on Learning Representations</i> .	1051
		1052
		1053
		1054
		1055
		1056
	Yuxun Tang, Yuning Wu, Jiatong Shi, and Qin Jin. 2024b. SingOMD: Singing oriented multi-resolution discrete representation construction from speech models . In <i>Interspeech 2024</i> , pages 2564–2568.	1057
		1058
		1059
		1060
	Jinchuan Tian, Chunlei Zhang, Jiatong Shi, Hao Zhang, Jianwei Yu, Shinji Watanabe, and Dong Yu. 2024. Preference alignment improves language model-based TTS. <i>arXiv preprint arXiv:2409.12403</i> .	1061
		1062
		1063
		1064
	Changhan Wang, Anne Wu, Jiatao Gu, and Juan Pino. 2021. CoVoST 2 and massively multilingual speech translation . In <i>Interspeech 2021</i> , pages 2247–2251.	1065
		1066
		1067
	Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang, Long Zhou, Shujie Liu, Zhuo Chen, Yanqing Liu, Huaming Wang, Jinyu Li, et al. 2023. Neural codec language models are zero-shot text to speech synthesizers. <i>arXiv preprint arXiv:2301.02111</i> .	1068
		1069
		1070
		1071
		1072
	Wei Wang, Shuo Ren, Yao Qian, Shujie Liu, Yu Shi, Yanmin Qian, and Michael Zeng. 2022. Optimizing alignment of speech and language latent spaces for end-to-end speech recognition and understanding. In <i>ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 7802–7806. IEEE.	1073
		1074
		1075
		1076
		1077
		1078
		1079
	Xiaofei Wang, Manthan Thakker, Zhuo Chen, Naoyuki Kanda, Sefik Emre Eskimez, Sanyuan Chen, Min Tang, Shujie Liu, Jinyu Li, and Takuya Yoshioka. 2024. Speechx: Neural codec language model as a versatile speech transformer . <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> .	1080
		1081
		1082
		1083
		1084
		1085
	Haibin Wu, Ho-Lam Chung, Yi-Cheng Lin, Yuan-Kuei Wu, Xuanjun Chen, Yu-Chi Pai, Hsiu-Hsuan Wang, Kai-Wei Chang, Alexander H. Liu, and Hung yi Lee. 2024a. Codec-SUPERB: An in-depth analysis of sound codec models . <i>Preprint</i> , arXiv:2402.13071.	1086
		1087
		1088
		1089
		1090
	Tongtong Wu, Linhao Luo, Yuan-Fang Li, Shirui Pan, Thuy-Trang Vu, and Gholamreza Haffari. 2024b. Continual learning for large language models: A survey . <i>Preprint</i> , arXiv:2402.01364.	1091
		1092
		1093
		1094
	Yihan Wu, Soumi Maiti, Yifan Peng, Wangyou Zhang, Chenda Li, Yuyue Wang, Xihua Wang, Shinji Watanabe, and Ruihua Song. 2024c. SpeechComposer: Unifying multiple speech tasks with prompt composition . <i>Preprint</i> , arXiv:2401.18045.	1095
		1096
		1097
		1098
		1099

1100	Zhanglin Wu, Jiabin Guo, Daimeng Wei, Zhiqiang Rao,	Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan,	1155
1101	Zongyao Li, Hengchao Shang, Yuanchang Luo, Shao-	Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. 2023b.	1156
1102	jun Li, and Hao Yang. 2024d. Improving the quality	SpeechGPT: Empowering large language models	1157
1103	of IWSLT 2024 cascade offline speech translation	with intrinsic cross-modal conversational abilities.	1158
1104	and speech-to-speech translation via translation hy-	In <i>Findings of the Association for Computational</i>	1159
1105	pothesis ensembling with nmt models and large lan-	<i>Linguistics: EMNLP 2023</i> , pages 15757–15773, Sin-	1160
1106	guage models. In <i>Proceedings of the 21st Interna-</i>	gapore. Association for Computational Linguistics.	1161
1107	<i>tional Conference on Spoken Language Translation</i>		
1108	<i>(IWSLT 2024)</i> , pages 46–52.		
1109	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng,	Xin Zhang, Dong Zhang, Shimin Li, Yaqian Zhou, and	1162
1110	Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan	Xipeng Qiu. 2024. Speechn tokenizer: Unified speech	1163
1111	Li, Dayiheng Liu, Fei Huang, et al. 2024a. Qwen2	tokenizer for speech language models. In <i>The Twelfth</i>	1164
1112	technical report. <i>arXiv preprint arXiv:2407.10671.</i>	<i>International Conference on Learning Representa-</i>	1165
		<i>tions.</i>	1166
1113	Dongchao Yang, Jinchuan Tian, Xu Tan, Rongjie Huang,	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,	1167
1114	Songxiang Liu, Haohan Guo, Xuankai Chang, Jia-	Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen	1168
1115	tong Shi, Jiang Bian, Zhou Zhao, et al. 2024b. Uni-	Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen	1169
1116	Audio: Towards universal audio generation with large	Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang,	1170
1117	language models. In <i>Forty-first International Confer-</i>	Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu,	1171
1118	<i>ence on Machine Learning.</i>	Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023.	1172
		A survey of large language models. <i>arXiv preprint</i>	1173
		<i>arXiv:2303.18223.</i>	1174
1119	Rong Ye, Chengqi Zhao, Tom Ko, Chutong Meng, Tao	Xinfa Zhu, Yuanjun Lv, Yi Lei, Tao Li, Wendi He,	1175
1120	Wang, Mingxuan Wang, and Jun Cao. 2023. Gigast:	Hongbin Zhou, Heng Lu, and Lei Xie. 2023. Vec-	1176
1121	A 10,000-hour pseudo speech translation corpus. In	Sok speech: Speech vectorization and tokeniza-	1177
1122	<i>INTERSPEECH 2023</i> , pages 2168–2172.	tion for neural speech generation. <i>arXiv preprint</i>	1178
		<i>arXiv:2310.07246.</i>	1179
1123	Neil Zeghidour, Alejandro Luebs, Ahmed Omran,	A Speech Tokenizer	1180
1124	Jan Skoglund, and Marco Tagliasacchi. 2021.		
1125	Soundstream: An end-to-end neural audio codec.	The tokenizer is optimized with 8-stream ($L = 8$)	1181
1126	<i>IEEE/ACM Transactions on Audio, Speech, and Lan-</i>	residual vector quantization (RVQ) layers, each	1182
1127	<i>guage Processing</i> , 30:495–507.	containing 1,024 tokens per codebook ($ \mathcal{B}_j =$	1183
1128	Heiga Zen, Viet Dang, Rob Clark, Yu Zhang, Ron J.	$1,024, (j = 1, \dots, L)$). The framerate is set to	1184
1129	Weiss, Ye Jia, Zhifeng Chen, and Yonghui Wu. 2019.	50Hz. Most hyperparameters related to the model	1185
1130	LibriTTS: A corpus derived from librispeech for text-	architecture are aligned with those in the original	1186
1131	to-speech. In <i>Interspeech 2019</i> , pages 1526–1530.	SoundStream paper (Zeghidour et al., 2021), while	1187
1132	Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing	the discriminator setups follow the DAC framework	1188
1133	Qu, Yong Jae Lee, and Yi Ma. 2023. Investigating the	(Kumar et al., 2023).	1189
1134	catastrophic forgetting in multimodal large language	We use a segment size of 24,000 samples (1.5	1190
1135	model fine-tuning. In <i>Conference on Parsimony and</i>	seconds) with a batch size of 6. The loss terms are	1191
1136	<i>Learning (Proceedings Track).</i>	consistent with those in the DAC paper. Both the	1192
1137	Jun Zhan, Junqi Dai, Jiasheng Ye, Yunhua Zhou,	generator and discriminator are optimized using the	1193
1138	Dong Zhang, Zhigeng Liu, Xin Zhang, Ruibin Yuan,	AdamW optimizer with a learning rate of 0.0002.	1194
1139	Ge Zhang, Linyang Li, et al. 2024. AnyGPT: Unified	We apply an exponential learning rate scheduler	1195
1140	multimodal LLM with discrete sequence modeling.	with a decay rate of 0.999.	1196
1141	<i>arXiv preprint arXiv:2402.12226.</i>	Due to the nature of RVQ, the initial streams	1197
1142	Binbin Zhang, Hang Lv, Pengcheng Guo, Qijie Shao,	in the codec typically carry the bulk of the sig-	1198
1143	Chao Yang, Lei Xie, Xin Xu, Hui Bu, Xiaoyu Chen,	nal information. To further enhance learning in	1199
1144	Chenchen Zeng, et al. 2022. Wenetspeech: A 10000+	these initial streams, we select the target bandwidth	1200
1145	hours multi-domain mandarin corpus for speech	from 0.5, 1, 1.5, 2.0, 4.0, sampled uniformly. Em-	1201
1146	recognition. In <i>ICASSP 2022-2022 IEEE Interna-</i>	pirically, this setup achieves better reconstruction	1202
1147	<i>tional Conference on Acoustics, Speech and Signal</i>	quality using only three codec levels ($L' = 3$), re-	1203
1148	<i>Processing (ICASSP)</i> , pages 6182–6186. IEEE.	ducing the modeling complexity in the codec-based	1204
1149	Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan,	speech LLM. We also introduce noise and rever-	1205
1150	Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. 2023a.	beration into 20% of the training data to improve	1206
1151	SpeechGPT: Empowering large language models	the model’s understanding capabilities. The signal-	1207
1152	with intrinsic cross-modal conversational abilities.	to-noise ratio (SNR) for these samples is randomly	1208
1153	In <i>Findings of the Association for Computational</i>		
1154	<i>Linguistics: EMNLP 2023</i> , pages 15757–15773.		

chosen from a range of 6.0 to 20.0 dB. The total training steps are 1.6M.

B Pre-training Dataset

B.1 Pre-training Data Details

The list of pre-training data is provided in Table 5.

B.2 Pre-training Data Simulation

Among recent speech LLMs, we identify two approaches to support multi-task training. One approach uses natural language prompts, enabling generalization across multiple tasks (Tang et al., 2024a; Gong et al., 2024; Chu et al., 2024a; Hu et al., 2024). These prompts can also serve as a bridge between the newly introduced speech modality and the original text, particularly when a pre-trained text LLM is involved.

The second approach uses a task template with either explicit task tokens or modality orders (Chu et al., 2023; Yang et al., 2024b; Wang et al., 2023; Chen et al., 2024b). While this method allows for more stable modeling across different tasks, it is less flexible. Most generation-oriented speech LLMs use this approach.

In our work, to best leverage the reasoning capabilities of the pre-trained text LLM, we opt for the natural language prompt method to connect different modalities.

To generate prompts with enough variations, we use OpenAI APIs for both ChatGPT powered by either GPT3.5 or GPT-4 (Achiam et al., 2023). The prompts are generated in both English and Mandarin. After the initial generation of 50 prompts for both English and Mandarin, we conduct manual filtering to remove unreasonable prompts, resulting in 25 English prompts and 25 Mandarin prompts per task. For the prompts used during pre-training, we limit the prompts to have a declarative format during initial generation, while we specify multiple formats for the prompts used during fine-tuning, including declaratives, interrogatives, and imperatives.

We define the pre-training data simulation policy as follows:

- For ASR, S2TT, TTS, and T2ST, the randomly selected natural language prompt can be in either language, regardless of the language used in the corresponding speech-text pair.
- For each sample in the TTS task, a portion of the target codec tokens is ran-

domly selected, with a duration ranging from $[\min(T_C/4, 2\text{seconds}), \min(T_C/2, 4\text{seconds})]$.

Additionally, we use ChatGPT to generate natural language prompts that specify the task of synthesizing speech in the speaking style of the target speech. These prompts are written in an imperative format, matching the language of the task-specific prompt, and are concatenated with the input text and acoustic conditioning in the following format: ‘<Text><Speaking Style Prompt><Acoustic Conditioning><Task Prompt>’.

C Model Hyper-parameters

Pre-training experiments are conducted with tensor parallelism set to 8, which enables the use of larger batch sizes during pre-training (Korthikanti et al., 2023). The gradient clip is set to 1.0. The normalization epsilon is set to $1e-5$. The global batch size is set to 640 with a sample sequence length of 4,096. The number of training steps is set to 40k. The model is trained on BF16. We use the distributed Adamw optimizer with a peak learning rate of $1e-5$ and a minimum learning rate of $1e-6$. The expanded vocabulary size is 155,012, considering padding tokens and 293 additional tokens for tensor shape adjustment to achieve tensor parallel training. The model has a parameter size of 943.5M.

Table 5: Continual pre-training dataset.

Dataset	Language	Data Type	Data Size (Hour)
Aishell{1-3} (Bu et al., 2017; Du et al., 2018; Shi et al., 2021)	ZH	Read	1,200
Wenetspeech (Zhang et al., 2022)	ZH	Various	10,000
Gigaspeech (Chen et al., 2021)	EN	Various	10,000
Librispeech (Panayotov et al., 2015)	EN	Read	1,000
MLS (Pratap et al., 2020)	EN	Read	44,000
TEDLIUM3 (Hernandez et al., 2018)	EN	Lecture	400
WSJ (Rottland et al., 1997)	EN	Read	140
Commonvoice (Ardila et al., 2020)	EN & ZH	Read	2,600
In-house	ZH	Various	70,000
Crawled Youtube subtitles	EN & ZH	-	-
Wikipedia	EN	-	-
The pile book corpus (Gao et al., 2020)	EN	-	-
In-house Mandarin data	ZH	-	-
In-house translation data	EN & ZH	-	-