

ELLA-V: Stable Neural Codec Language Modeling with Alignment-Guided Sequence Reordering

Yakun Song¹, Zhuo Chen², Xiaofei Wang³, Ziyang Ma¹, Xie Chen^{1*}

¹Shanghai Jiao Tong University, Shanghai, China

²ByteDance Inc., USA

³Microsoft, One Microsoft Way, Redmond, USA
{ereboas, chenxie95}@sjtu.edu.cn

Abstract

The language model (LM) approach based on acoustic and linguistic prompts, such as VALL-E, has achieved remarkable progress in the field of zero-shot audio generation. However, existing methods still have some limitations: 1) repetitions, transpositions, and omissions in the output synthesized speech due to limited alignment constraints between audio and phoneme tokens; 2) challenges of fine-grained control over the synthesized speech with autoregressive (AR) language model; 3) infinite silence generation due to the nature of AR-based decoding, especially under the greedy strategy. To alleviate these issues, we propose ELLA-V, a simple but efficient LM-based zero-shot text-to-speech (TTS) framework, which enables fine-grained control over synthesized audio at the phoneme level. The key to ELLA-V is interleaving sequences of acoustic and phoneme tokens, where phoneme tokens appear ahead of the corresponding acoustic tokens. The experimental findings reveal that our model outperforms baselines in terms of accuracy and delivers more stable results using both greedy and sampling-based decoding strategies.

Demo & Code — <https://ereboas.github.io/ELLAV/>

1 Introduction

Recently, deep generative AI has achieved remarkable results in various tasks, leading to the emergence of many transformative real-world applications (Brown et al. 2020; Ramesh et al. 2022; Ho, Jain, and Abbeel 2020). There have been rapid developments in the field of speech synthesis as well. In particular, zero-shot TTS technology has gained increasing attention because it can synthesize high-quality target voices without the need of specified speaker’s training data. As a state-of-the-art generative model family, diffusion models (Sohl-Dickstein et al. 2015; Ho, Jain, and Abbeel 2020; Song and Ermon 2020) progressively add noise to the training data and then learn the reverse process to generate samples. By leveraging diffusion models and their variants (Sohl-Dickstein et al. 2015; Ho, Jain, and Abbeel 2020; Lipman et al. 2023), many works have successfully applied them to the audio domain (Shen et al. 2023;

Le et al. 2023; Kim et al. 2024; Vyas et al. 2023). However, they may also face the limitation in training an additional duration predictor. Another major class of generative models is language modeling based on Transformer (Vaswani et al. 2017a). AR language models use a decoder-only architecture to predict the next token in a sequence as the training objective, which has demonstrated extremely powerful few-shot and zero-shot capabilities in many generative tasks (Brown et al. 2020; Thoppilan et al. 2022; Chowdhery et al. 2023). In light of this, VALL-E (Wang et al. 2023a) and subsequent works (Kharitonov et al. 2023; Rubenstein et al. 2023; Wang et al. 2023b) have successfully employed decoder-only language model for zero-shot TTS. These approaches first quantize the speech signal into a series of discrete acoustic tokens. Subsequently, they employ an AR language model to predict coarse-grained acoustic tokens, eliminating the necessity for explicit duration predictors or speaker encoders. Once trained on a large-scale corpus, these approaches are capable of synthesizing speech with competitive fidelity and naturalness in a zero-shot manner.

While VALL-E and its variants have achieved numerous impressive milestones, they still possess certain limitations that impact practical deployment. For instance, existing methods (Wang et al. 2023a; Kharitonov et al. 2023) directly concatenate phoneme tokens and acoustic tokens as a whole sequence to train language models. In this way, the alignment between audio and phoneme sequences is completely learned through the self-attention in the transformer, making it potentially unstable as self-attention does not explicitly capture the monotonic alignment between audio and phoneme. Additionally, the decoder-only language model architecture can lead to potential attention degradation issues (Fu et al. 2023), where the alignment quality between the target audio sequence and the source phoneme sequence deteriorates as the generated sequence increases, resulting in inaccurate or low-quality speech outputs.

Another limitation stems from the nature of AR language modeling. Specifically, given a sequence \mathbf{x} , the standard AR language model factorizes the likelihood $p(\mathbf{x})$ over the dimensions of \mathbf{x} via the chain rule $p(\mathbf{x}) = \prod_{t=0}^T p(x_t | \mathbf{x}_{<t})$. AR models predict the current tokens solely based on the historical tokens without users’ control in the inference process, and sometimes generate semantic repetitions or incoherence in the generated output (Yang et al. 2019; Brown

*Corresponding author.

Top- p	WER%	INF%
1	5.47	0.00
0.99	5.00	0.20
0.95	10.99	19.06
0.9	20.85	41.43
0.7	37.71	76.76
0.4	46.59	84.39
0.0 (greedy)	49.26	87.29

Table 1: Comparison of VALL-E’s zero-shot TTS performance across various top- p thresholds in nuclear sampling. INF% denotes the probability of *infinite silence*, which refers to instances where generation continues without stopping when its duration exceeds twice the original length.

et al. 2020). In the TTS task, correspondingly, VALL-E cannot directly determine which segment of the output audio corresponds to which prompt phoneme, thus there is no trivial way to promptly detect and prevent issues occurring in the generation process. These drawbacks can manifest as meaningless phoneme repetitions, transpositions, omissions, or even catastrophic *infinite silence*, i.e., during the process of generation, the model anomalously outputs silence or noise tokens for an extended period of time without stopping. Specifically, Table 1 demonstrates the word error rate (WER) and the probability of the *infinite silence* in VALL-E samples at different threshold top- p for nuclear sampling (Holtzman et al. 2019). The detailed experimental setup is described in Section 4. Notably, a shift in the decoding strategy of VALL-E from fully sampling-based to fully greedy-based leads to a marked decline in sample quality. It should be emphasized that while sampling-based stochastic decoding strategies have advantages in terms of synthesis diversity, deterministic decoding strategies (e.g., beam search and its variants) are more suitable for cases where there is less tolerance for synthesis errors and more emphasis on fluency and coherence (Ippolito et al. 2019).

Faced with the pros and cons of the existing methods, we introduce ELLA-V, a simple but effective language model approach for zero-shot TTS. ELLA-V proposes a generalized AR (GAR) language model to generate the first layer of residual vector quantizer (RVQ) codes of a neural codec model. Then as with VALL-E, ELLA-V employs a non-autoregressive (NAR) language model to obtain codes of the other RVQs. Our core innovation lies in 3 fold:

- Firstly, ELLA-V inserts phone tokens into the corresponding positions of the acoustic sequence. Unlike existing methods, Connecting phoneme tokens with their corresponding acoustic tokens can help the language model capture the alignment between phoneme and acoustic modalities in local dependencies.
- Secondly, instead of maximizing the expected log-likelihood of the hybrid sequence under a conventional casual mask, ELLA-V computes loss only on acoustic tokens and special tokens. This training objective provides a natural way to have fine-grained control in inference: ELLA-V’s GAR model always maintains aware-

ness of the phoneme it is currently synthesizing, allowing it to promptly detect and truncate any abnormal phoneme to avoid any possible *infinite silence* issue.

- Thirdly, we further propose an improvement to the input sequence. We introduce *local advance*, which involves shifting the EOP token and the next-word phoneme token a few frames ahead. By advancing these special tokens, the GAR model can better utilize local dependencies to predict the pronunciation of the current phoneme.

Experimental results, using comparable model configurations and 960 hours of speech data from LibriSpeech (Panayotov et al. 2015) as a training set, demonstrate the superiority of ELLA-V. Compared to the cutting-edge zero-shot TTS baseline systems, ELLA-V significantly improves the accuracy of synthesized speech, and demonstrates comparable or superior speaker similarity and speech naturalness on a series of subjective and objective experiments. Notably, ELLA-V works well on a wide spectrum of decoding strategies – even greedy decoding, and still has a substantially better speech accuracy than the best of VALL-E. We further conducted ablation experiments to investigate the effects of our proposed modifications. The results indicate that the *global advance* in ELLA-V significantly improves the model’s performance, while the local advance enhances the stability of the generated output.

2 Related Work

speech synthesis Speech synthesis has long been a significant topic in the fields of natural language processing, and speech processing. Early methods were based on Statistical Parametric Speech Synthesis (Zen, Tokuda, and Black 2009), typically involving complex components such as text analysis models, acoustic models, and vocoders. Later, end-to-end neural TTS models were introduced, which synthesize Mel spectrograms and employ a vocoder (Oord et al. 2017) for speech synthesis (Wang et al. 2017). Some methods, utilizing techniques such as VAE (Hsu et al. 2019; Lee, Shin, and Jung 2022), flow (Miao et al. 2020; Kim et al. 2020), diffusion (Jeong et al. 2021; Vyas et al. 2023), and others (Wu and Shi 2022), have achieved promising performance in end-to-end speech synthesis, although explicit duration predictors or speaker embedders are usually necessary. On the other hand, models like VALL-E (Wang et al. 2023a) and AudioLM (Borsos et al. 2023a) utilize autoregressive Transformers to model discrete audio tokens, achieving great in-context learning performance. When it comes to zero-shot speech synthesis, autoregressive Transformer-based models can predict and generate audio without the need for an additional duration model, which strikes a favorable balance between efficiency and performance, and has been garnering increasing attention.

3 Method

3.1 Overview

Fig. 1 demonstrates the overall architecture of ELLA-V. ELLA-V primarily follows a two-stage framework similar to VALL-E, considering zero-shot TTS as a conditional

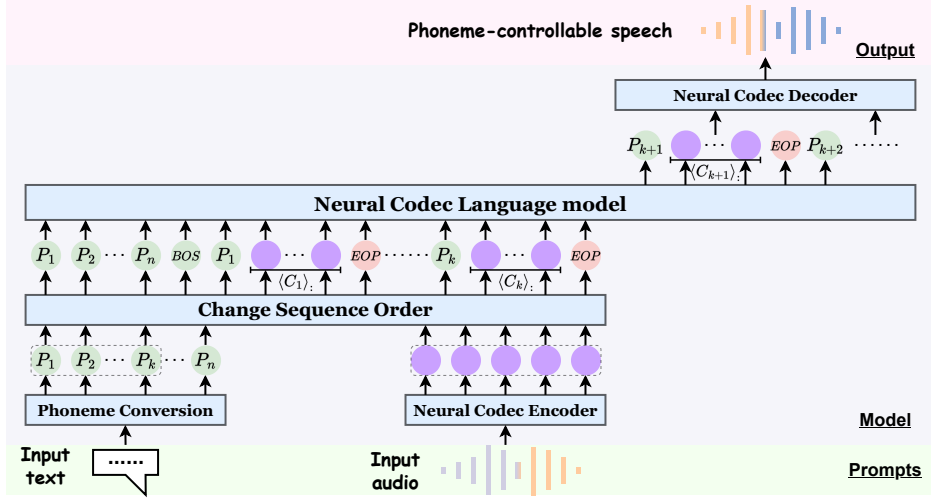


Figure 1: The overall architecture of ELLA-V. Input an audio prompts and text prompts, ELLA-V first changes sequence order – sandwiching each phoneme’s audio $\langle C_k \rangle$, between the k -th phoneme and a EOP token and prepending the phoneme sequence to the beginning. By learning on the mixed sequence, ELLA-V can generate audio sequence of the text prompts while maintaining the acoustic and environmental conditions of the audio prompts.

codec language modeling task. ELLA-V maps input text prompts and speech prompts into a unified vocabulary space with a text encoder and a neural codec, respectively. Different from VALL-E, an additional sequence order rearranging step is performed to the text-audio token sequence, after which, ELLA-V utilizes a decoder-only language model to learn to perform conditional generation on the hybrid sequences of phoneme and audio tokens. Detailed information about the language model will be presented in Section 3.2.

To obtain discrete audio representations, we employ a pre-trained neural audio codec model, EnCodec (Défossez et al. 2023), following VALL-E (Wang et al. 2023a). EnCodec transforms 24 kHz raw waveforms into 75 Hz discrete tokens using L RVQ layers. In our experiments, we use the same settings as VALL-E, with $L = 8$ and the codebook size is 1024. In this setting, each second of the waveform is represented by 75×8 discrete tokens from RVQ.

To obtain phoneme sequences, we apply the Montreal Forced Aligner (MFA) (McAuliffe et al. 2017) to the input audio. Notably, MFA not only serves as a text tokenizer but also extracts alignment relationships between phonemes and the corresponding speech. The forced alignment information is essential for ELLA-V to **change sequence order**. In Section 3.2, we will provide a detailed explanation of how this information is used to construct the target sequence.

3.2 Training: Codec Language Model

ELLA-V employs a Generalized Autoregressive Codec language model for the prediction of the first quantization layer in the EnCodec, which corresponds to capturing semantic information and coarse-grained acoustic profiles. Subsequently, a non-autoregressive language model is utilized to generate codes for the subsequent quantization layers, aimed at reconstructing fine-grained acoustic details. Specifically, given a speech corpus $\mathcal{D} = \{\mathbf{x}_i, \mathbf{y}_i\}$, where \mathbf{x}_i represents the

i -th audio sample, and \mathbf{y}_i is its text transcription. We utilize the EnCodec to extract the discrete representation of \mathbf{x} , formulated as $\mathbf{C}^{T \times 8} = \text{EnCodec}(\mathbf{x})$, where \mathbf{C} represents the two-dimensional acoustic code matrix, and T is the down-sampled utterance length.

We employ MFA to obtain the phoneme sequence $\mathbf{P}_{1:n}$ corresponding to the transcription \mathbf{y} , while also extracting forced alignment information between the audio \mathbf{x} and the transcription \mathbf{y} : $(\mathbf{P}_{1:n}, \mathbf{l}_{1:n}) = \text{MFA}(\mathbf{x}, \mathbf{y})$, where n is the number of phonemes of the audio sample \mathbf{x} , and l_i denotes the length of the i -th phoneme of the discrete audio sequence. MFA treats silence also as a kind of phoneme, so that the original audio sequence is partitioned into n consecutive intervals corresponding to n phonemes. Specifically, let $\langle \mathbf{C}_i \rangle^{l_i \times 8}$ represent the audio sequence corresponding to the i -th phoneme, \mathbf{C} is the concatenation of $\langle \mathbf{C}_i \rangle$, and we have $\langle \mathbf{C}_k \rangle_{1:l_k} = \mathbf{C}_{\sum_{i=1}^{k-1} l_i + 1 : \sum_{i=1}^k l_i}$

After quantization, we utilize the EnCodec decoder to reconstruct the audio waveform from the discrete acoustic sequence \mathbf{C} , formulated as $\hat{\mathbf{x}} \approx \text{DeCodec}(\mathbf{C})$.

For the zero-shot TTS task, the optimization objective is $\max p(\mathbf{C} | \mathbf{P}, \hat{\mathbf{C}})$, where $\hat{\mathbf{C}}$ is the acoustic prompt of the unseen speaker. We use language modeling to generate acoustic tokens, by learning on the mixed sequence composed of phonemes and codec codes, consistent with previous works (Wang et al. 2023a; Rubenstein et al. 2023).

Unlike existing approaches, ELLA-V does not concatenate phoneme tokens and acoustic tokens directly to form the target sequence for training the language model. Instead, ELLA-V **interleaves phoneme and acoustic tokens** in order to make it easier for language models to learn the alignment between audio and text. Specifically, we insert each phoneme token P_i (except the silence phoneme) into the corresponding position of the audio sequence, so that each

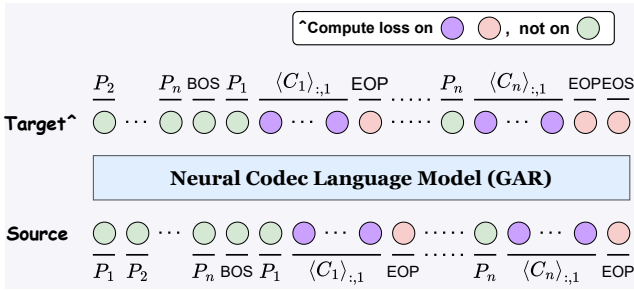


Figure 2: The illustration of Generalized Autoregressive language model of ELLA-V.

phoneme’s audio $\langle C_i \rangle$ is sandwiched between P_i and EOP tokens. We also prepend the phoneme sequence to the beginning of the mixed sequence, which is referred to as *global advance*. In Section 3.4, we further propose a variant sequence order with higher generation stability, named *local advance*, which moves the non-acoustic tokens of the sequence several frames forward.

Generalized Autoregressive (GAR) Codec Language Model As shown in Figure 2, ELLA-V first constructs a hybrid sequence $\mathbf{H}_{:,1}$ of acoustic and phoneme tokens, structured as: $[P_1, P_2, \dots, P_n, \text{BOS}, P_1, \langle C_1 \rangle_{:,1}, \text{EOP}, \dots, P_n, \langle C_n \rangle_{:,1}, \text{EOP}, \text{EOS}]$. It is worth noting that the MFA (Montreal Forced Aligner) treats silence as a distinct phoneme, whereas our phoneme sequence \mathbf{P} exclusively comprises phonemes other than silence. To clarify, we retain the acoustic component associated with silence but do not sandwich it with an EOP and a specific silence phoneme, nor do we use a silence phoneme in the *global advance* part.

We design a GAR language model to learn the continuation task on the hybrid sequence, to generate the discrete acoustic code sequence $\mathbf{C}_{:,1}$. The GAR model consists of multiple Transformer decoder layers (Vaswani et al. 2017b). After training, it can generate discrete audio codes for a specified text prompt and acoustic prompt. GAR is also responsible for predicting EOP and EOS to indicate the conclusion of a phoneme and the entire sentence, respectively.

The optimization of GAR is achieved by maximizing the likelihood of the acoustic part $\mathbf{C}_{:,1}$ of the hybrid sequence $\mathbf{H}_{:,1}$, as well as the special EOP and EOS tokens. Under forward factorization, this process is formulated as:

$$\begin{aligned}
& \max_{\theta_{GAR}} \log p(\tilde{\mathbf{C}}_{:,1} | \mathbf{P}; \theta_{GAR}) \\
&= \sum_{i=1}^n \sum_{t=0}^{l_i} \log p(\langle \tilde{\mathbf{C}}_i \rangle_{t,1} | \langle \tilde{\mathbf{C}}_i \rangle_{<t,1}, \langle \tilde{\mathbf{C}}_{<i} \rangle_{:,1}, \\
& \quad \mathbf{P}; \theta_{GAR}) \tag{1} \\
&= \sum_{\substack{t=0 \\ \mathbf{H}_{t,1} \neq \text{BOS} \\ \mathbf{H}_{t,1} \notin \{\mathbf{P}\}}}^{T_H} \log p(\mathbf{H}_{t,1} | \mathbf{H}_{<t,1}; \theta_{GAR})
\end{aligned}$$

where \mathbf{H} has a size of $T_H \times 8$, $\{\mathbf{P}\}$ denotes the phoneme set, $\langle \tilde{\mathbf{C}}_i \rangle$ is the concatenation of $\langle C_i \rangle$ along with its broadcast

trailing EOP and/or EOS tokens, $\tilde{\mathbf{C}}$ is then the concatenation of $\langle C_i \rangle$, and θ_{GAR} represents neural network parameters of GAR model. The factorization of the training objective naturally encapsulates the core intuition of the GAR model: GAR generates the audio sequence phoneme-by-phoneme. GAR produces maximum likelihood predictions for each phoneme token successively, indicating the end of generating a specified phoneme by predicting EOP. Through *global advancement*, GAR can directly infer the next phoneme to be generated without relying on network predictions. After the prediction for the last phoneme is completed, GAR stops the generation process by predicting EOS. The generated sequence by GAR is **self-aligned**, as it can instantly know the corresponding position of any generated acoustic token in relation to the phoneme prompt.

During training, we apply a bidirectional mask to the phoneme sequence before the BOS in the hybrid sequence, while a unidirectional mask is used for the part after BOS. We frame the training as a next-token-prediction language modeling task on the hybrid sequence. However, it’s important to note that the model does not predict phonemes (or BOS). As shown in Figure 2, we only compute loss when the token to be predicted is not a phoneme (or BOS). During inference, after the model predicts an EOP for a phoneme, the next phoneme token is directly appended to the end of the sequence, which will be further discussed in Section 4.

Non-Autoregressive (NAR) Codec Language Model In the second stage, the NAR language model is employed to predict the codes from the second to the last quantization layers in parallel. The input-output sequence construction of the NAR model follows the same pattern as used in the GAR model. Specifically, the i -th column $\mathbf{H}_{:,i}$ of the hybrid sequence matrix \mathbf{H} is structured as: $[P_1, P_2, \dots, P_n, \text{BOS}, P_1, \langle C_1 \rangle_{:,i}, \text{EOP}, \dots, P_n, \langle C_n \rangle_{:,i}, \text{EOP}, \text{EOS}]$.

And in practice if P_i represents the silence, $\mathbf{C}_{:,i}$ will not be sandwiched by P_i and EOP.

The NAR model takes the previously generated hybrid sequence of the previous $j - 1$ layers as input and predicts the codes of the j -th layer in parallel, formulated as:

$$\begin{aligned}
& \max_{\theta_{NAR}} \sum_{j=2}^8 \log p(\mathbf{C}_{:,j} | \mathbf{H}_{:,<j}, \mathbf{P}; \theta_{NAR}) \\
&= \sum_{j=2}^8 \sum_{\substack{t=0 \\ \mathbf{H}_{t,j} \in \{\mathbf{C}_{:,j}\}}}^{T_H} \log p(\mathbf{H}_{t,j} | \mathbf{H}_{:,<j}, \mathbf{P}; \theta_{NAR}) \tag{2}
\end{aligned}$$

where $\{\mathbf{C}_{:,j}\}$ denotes the acoustic token set of the j -th quantizer. The embeddings of tokens from the previous $j - 1$ quantizers are summed up to feed the NAR model to predict the j -th layer.

3.3 Inference

ELLA-V can use a short clip of speech from an unseen speaker as an acoustic prompt to synthesize speech for a specified text prompt. Figure 3 illustrates the inference process of the GAR model. While VALL-E may get stuck in an infinite loop during inference, resulting in the synthesis of

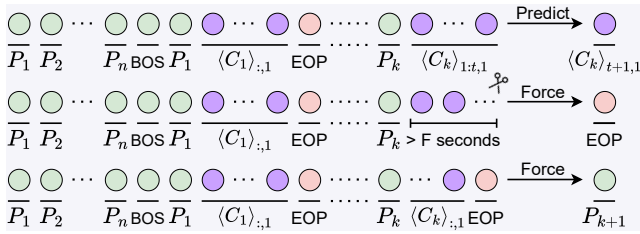


Figure 3: Illustration of the inference process of ELLA-V.

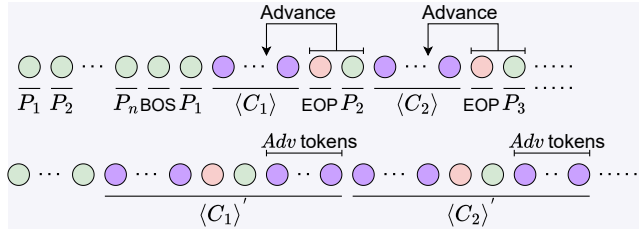


Figure 4: *Local advance*. A phoneme can locally have access to information about the next phoneme token advanced by *Adv* frames, allowing it to anticipate the upcoming phoneme token’s characteristics.

either *infinite silence* or repetitive pronunciation, ELLA-V is capable of generating EOP and promptly truncating abnormally long phonemes. Following an EOP, we can directly append the next phoneme token to the end of the generated sequence, ensuring the proper generation of speech without abnormal pauses or repetitions. For the GAR model, we employ a sampling-based decoding strategy, whereas for the NAR model, we use a greedy decoding approach to strike a balance between efficiency and performance.

3.4 Local Advance

One intuition is that the pronunciation of a phoneme is strongly related to the phonemes around it. However, due to the autoregressive nature of the GAR model, an acoustic token cannot attend to the following phoneme tokens, even though we can leverage the transformer’s ability to model long-term dependencies through *global advance* to provide complete context for the acoustic token generation. To further harness the powerful capability of the transformer in modeling local dependencies, ELLA-V introduces an additional change in the sequence order based on Section 3.2. Specifically, we move the phoneme token and the EOP token ahead by a few frames, referred to as *local advance*.

4 Experiments

4.1 Experimental Setup

Data & Tasks: We trained ELLA-V using the publicly available Librispeech (Panayotov et al. 2015) 960h training dataset. We utilized Montreal Forced Aligner (MFA)¹ (McAuliffe et al. 2017) to obtain forced alignment

¹<https://github.com/MontrealCorpusTools/Montreal-Forced-Aligner>

information for the audio-transcription pairs. Sentences with unrecognized or unknown phones by MFA were excluded. The open-source 24kHz checkpoint² of EnCodec(Défossez et al. 2023) was used as the codec to generate discrete acoustic tokens. The LibriSpeech training data was upsampled to 24 kHz before feeding it into EnCodec.

In evaluating the model, two zero-shot TTS tasks were considered. For the zero-shot TTS continuation task, we adhered to methodologies established by previous works (Wang et al. 2023a,c), selecting examples ranging from 4 seconds to 10 seconds from the LibriSpeech test-clean dataset as our test set. In this task, we used the complete phoneme transcription as the text prompt and the first 3 seconds of the test audio sample as the acoustic prompt. The model was required to generate continuations.

For the zero-shot TTS cross-speaker task, we designed a hard case set comprising 100 hard sentences. They included challenging phonetic patterns, alliteration, and unusual (abnormal) combinations of words that might pose difficulties for a TTS system to generate natural-sounding speech. In this case, we randomly picked 3-second sentences from the LibriSpeech test-clean subset as the acoustic prompt. We then concatenated the transcription of this segment and the target phoneme sequence in the hard case set to form the text prompt. The model was tasked with cloning the voice of the speaker to say the specified target text in the hard case set.

Training Configuration: For both GAR and NAR models, we stacked 12 Transformer decoder layers with an embedding dimension of 1024, a hidden state dimension of 1024, and a feed-forward layer dimension of 4096. All models were trained in parallel using 8 NVIDIA Tesla V100 GPUs with a batch size of 16384 tokens for GAR and 12288 tokens for NAR per GPU, respectively, learning a total of 320k steps. We used the AdamW optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-9}$. We employed an inverse-sqrt learning rate scheduler with warm-up. For the first 32000 updates, we linearly increased the learning rate from 10^{-7} to a peak of 5×10^{-4} . The weight decay was 0.01.

Baseline: In our research, we benchmarked the performance of zero-shot speech synthesis against VALL-E (Wang et al. 2023a). VALL-E was originally trained on a substantial 60k hours of audio from the Librilight dataset (Kahn et al. 2020). To ensure a rigorous evaluation, we reproduced the VALL-E model and adapted it to train on the LibriSpeech 960h dataset. We also adjusted the model dimensions and the number of layers to match the parameter settings of ELLA-V and VALL-E. Both GAR (or AR) and NAR models of VALL-E and ELLA-V have 154.3M parameters. Moreover, to mitigate any potential bias introduced by the audio codec, we pre-processed the authentic speech samples using EnCodec’s encoder and decoder. We include the result for EnCodec reconstructed speech for reference, denoted as Ground-Truth EnCodec. Two cutting-edge Non-AR-based TTS models were also included for comparison with our model: YourTTS (Casanova et al. 2022) trained on the 585-hour LibriSpeech subset LibriTTS (Zen et al. 2019), and

²<https://github.com/facebookresearch/encodec>

Models	WER(%) (↓)	SPK (↑)	CMOS	SMOS	Training Set	ASR Model
Ground Truth	1.41	0.923	0.29	4.39	/	C-T
Ground Truth-Encoder [†]	1.62	0.913	0.22	4.33	/	C-T
SoundStorm	2.99	/	/	/	LibriLight (60k hours)	C-T
YourTTS	6.34	0.822	/	/	LibriTTS (585 hours)	C-T
VALL-E	5.00	0.868	0.00	3.56	LibriSpeech (960 hours)	C-T
ELLA-V(<i>ours</i>)	2.28	0.870	0.10	3.56	LibriSpeech (960 hours)	C-T

Table 2: Main results of ELLA-V on zero-shot TTS continuation task. [†] indicates that ground-truth audios were passed through the encoder and decoder of Encoder to evaluate the influence of neural audio codec. The ASR model was used to obtain transcriptions of synthetic speech to calculate the WER. C-T denotes Conformer-Transducer. Since it is not open-sourced and has the same evaluation settings as this article, SoundStorm’s WER was derived directly from the original text.

Models	WER(%)	Sub(%)	Del(%)	Ins(%)
VALL-E	28.39	17.79	5.36	5.24
YourTTS	17.61	11.00	5.04	1.57
ELLA-V	12.79	7.76	3.40	1.63

Table 3: WER comparison on 100 particularly hard synthesis cases. Sub, Del, and Ins refer to Substitution, Deletion, and Insertion error rates, respectively.

Models	WER(%) (↓)	SPK (↑)
VALL-E	5.00	0.868
ELLA-V	2.28	0.870
ELLA-V-noglobal	5.00	0.859
ELLA-V-nophn	3.51	0.868

Table 4: The ablation study to investigate the impact of global and local phoneme information.

SoundStorm (Borsos et al. 2023b) trained on LibriLight.

Evaluation Metrics: We evaluated our system with several objective metrics. Speaker similarity (SPK) and WER served as our primary measures. SPK was assessed using the fine-tuned WavLM-TDNN model³ (Chen et al. 2022), scoring similarity on a scale of -1 to 1, with values above 0.86 indicate the same speaker identity (This value comes from the release model card page). The WER was determined by comparing the synthesized speech to the original text using the Conformer-Transducer model⁴ (Gulati et al. 2020).

In addition to these standard metrics, we introduced two novel measures: INF% and CUT%. INF% quantified the frequency of generating infinitely long audio, indicative of a failure in synthesis. It is used to measure the likelihood of the model falling into abnormal repetition (such as infinite silence). A higher INF% indicates poorer stability in the generated output of the model. In the practical implementation, INF% referred to the proportion of sentences for which generation was not stopped when the length of the generated audio reached twice the original, serving as a proxy for infinite generation. On the other hand, as discussed in the previous session, the design of ELLA-V enables control of the duration for each phoneme during inference, thus avoiding the synthesis failure. In our experiments, we forcibly truncate the synthesis of phonemes with a length greater than 0.4 seconds. CUT% is used to measure the frequency of forced cuts of phonemes in synthesis by ELLA-V. For each objective metric, we reported average values over three experimental

³<https://huggingface.co/microsoft/wavlm-base-plus-sv>

⁴https://huggingface.co/nvidia/stt_en_conformer_transducer_xlarge

runs with different random seeds.

For subjective analysis, we relied on the mean opinion score (MOS). 30 test samples were chosen for this purpose, with each sample being evaluated by at least 15 listeners for aspects like naturalness and speaker similarity. The comparative mean opinion score (CMOS) and the similarity mean opinion score (SMOS) were the key subjective metrics used. SMOS was rated on a 1 to 5 scale, in 0.5-point increments, to gauge speaker similarity, while CMOS, ranging from -1 to 1, assessed the overall naturalness and quality of the synthesized speech against the baseline.

4.2 Results

Zero-Shot TTS Continuation Task. We present the evaluation results in Table 2. First, regarding speaker similarity, both subjective (SMOS) and objective (SPK) results revealed that ELLA-V and VALL-E performed similarly and surpassed baseline Non-AR-based methods, which can be attributed to their shared backbone approach, combining (G)AR and NAR. Meanwhile, CMOS testing shows that ELLA-V achieved a +0.10 score, demonstrating a higher generation quality (i.e., naturalness) compared to VALL-E. Additionally, WERs calculated between the recognized text of synthesized audio and the ground-truth text show that ELLA-V is significantly better than VALL-E (2.28 versus 5.00), and better than SoundStorm (2.28 versus 2.99) trained on the dataset over 60 times larger. This underscores ELLA-V’s enhanced capability in synthesizing higher-quality and more robust speech. Overall, ELLA-V substantially improved the synthesis accuracy and robustness of the language model-based TTS framework without affecting the naturalness and speaker similarity. This conclusion is not

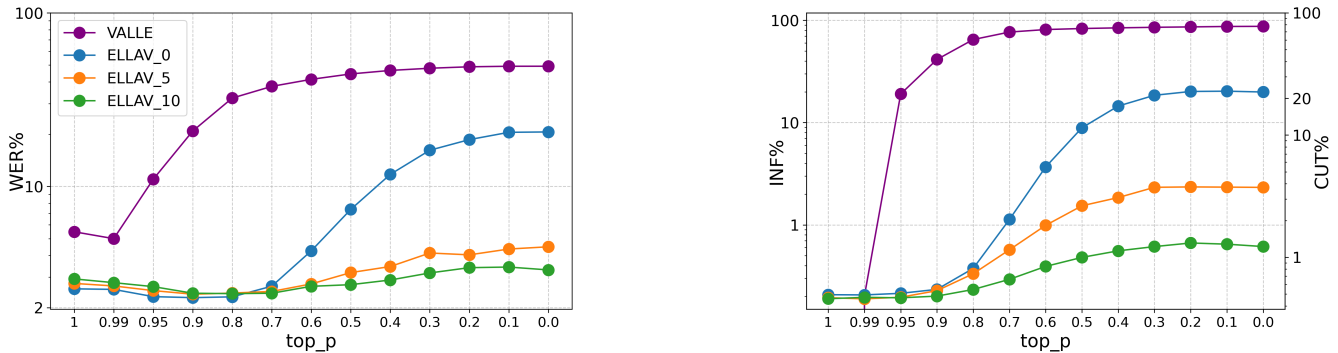


Figure 5: Ablations on decoding strategies. The left figure demonstrates the trends of WER for VALL-E, and ELLA-V with different local advance, with respect to the variations in top_p in nuclear sampling. The right figure shows trends of INF for VALL-E, and CUT for ELLA-V with different local advance, with respect to the variations in top_p in nuclear sampling.

only corroborated by this easy continuation task, but also validated via the challenging sets in the subsequent section.

Zero-Shot TTS Cross-Speaker Task on Hard Cases.

VALL-E utilized a traditional AR model that frequently resulted in alignment errors, including repetitions, transpositions, and omissions, particularly in more challenging synthesis cases (details of the challenging synthesis set can be found in supplementary material). Table 3 presents the WER results on the 100 particularly hard synthesis sentences. In contrast to the baselines, ELLA-V demonstrates markedly lower WER, signifying its enhanced robustness. This substantial reduction in errors translates to more accurate and reliable voice synthesis applications, significantly improving user experience in real-world scenarios.

Regarding VALL-E’s tendency to fall into infinite silence, an intuitive explanation is that the silence patterns in the training data are relatively simple and many of them are repetitive. In this case, a traditional language model is prone to overfitting to these patterns. During testing, when the model encounters silence, it assigns a high probability to silence. This leads to issues such as beam search, which is based on maximum likelihood, getting stuck in a loop. However, ELLA-V does not face this problem.

Analysis of Decoding Strategies and Local Advance. To demonstrate the stability of ELLA-V under different decoding strategies, we conducted an ablation study, testing the decoding performance with different top-p values for nuclear sampling, by varying p. The experiments were performed under local advance values of 0, 5, and 10. The results are shown in Figure 5. We can observe that as top_p decreases, the accuracy of VALL-E’s synthesized speech significantly decreases. At this point, VALL-E is more prone to generating a large number of overfit silence tokens, leading to a significant increase in INF%. And compared to VALL-E, the audio synthesized by ELLA-V is less sensitive to rate changes in the top_p sampling strategy, whose WER consistently outperforms VALL-E. When the local advance is set to 5 or 10 tokens, the generated audio exhibits significant stronger robustness. On the other hand, as shown in Figure

5 (right), as top_p decreases, VALL-E tends to get stuck in infinite loops of failed generation, while the generation of ELLA-V remains significantly stable. Moreover, ELLA-V can promptly handle (truncate) the synthesis of exceptional phonemes, resulting in significantly higher robustness.

Ablation Study. In this paragraph, we conduct ablation experiments. (1) To investigate the impact of global phoneme information on synthesized speech, we removed the global phoneme sequence at the beginning of the trained sequence (abbr. ELLA-V-noglobal). (2) To investigate whether it is necessary to provide the specific phoneme token before its corresponding acoustic tokens during both training and inference, rather than just using the EOP separator, we removed all phoneme tokens following BOS in the mixed sequence (abbr. ELLA-V-nophn). The experimental results are shown in Table 4. It is observed that the accuracy of synthesized speech significantly deteriorated either when global phoneme tokens were not used or when local phoneme tokens were disabled within the hybrid sequence. It is also notable that even in the absence of global advance (i.e., in the ELLA-V-noglobal configuration), the SPK and WER of the synthesized audio were comparable to those of VALL-E. These findings indicate the importance of both local and global information in achieving more accurate synthesized audios, meanwhile, combining both of them potentially leads to further enhancements in accuracy.

5 Conclusion

In this paper, we introduce ELLA-V, a simple and efficient two-stage zero-shot TTS framework based on language modeling. By learning interleaved sequences of acoustic and text tokens, our proposed GAR model can provide fine-grained control over synthesized audio at the phoneme level and can better leverage local dependencies to predict the pronunciation of the current phoneme. Experimental results demonstrate that ELLA-V achieves higher accuracy and more stable results under different threshold top-p for nuclear sampling. We aspire for this work to advance research in enhancing the robustness of speech generation.

Acknowledgments

This work was supported by the National Natural Science Foundation of China (No. 62206171 and No. U23B2018), Shanghai Municipal Science and Technology Major Project under Grant 2021SHZDZX0102, and the Science and Technology Innovation (STI) 2030-Major Projects under Grant 2022ZD0208700.

References

- Borsos, Z.; Marinier, R.; Vincent, D.; Kharitonov, E.; et al. 2023a. AudioLM: A language modeling approach to audio generation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*.
- Borsos, Z.; Sharifi, M.; Vincent, D.; Kharitonov, E.; et al. 2023b. Soundstorm: Efficient parallel audio generation. *arXiv preprint arXiv:2305.09636*.
- Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; et al. 2020. Language Models are Few-Shot Learners. In Larochelle, H.; Ranzato, M.; Hadsell, R.; Balcan, M.; and Lin, H., eds., *Proc. NeurIPS*.
- Casanova, E.; Weber, J.; Shulby, C. D.; Junior, A.; et al. 2022. Yourtts: Towards zero-shot multi-speaker tts and zero-shot voice conversion for everyone. In *Proc. ICML*.
- Chen, S.; Wang, C.; Chen, Z.; Wu, Y.; et al. 2022. WavLM: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing*.
- Chowdhery, A.; Narang, S.; Devlin, J.; Bosma, M.; et al. 2023. PaLM: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240): 1–113.
- Défosses, A.; Copet, J.; Synnaeve, G.; and Adi, Y. 2023. High Fidelity Neural Audio Compression. *Transactions on Machine Learning Research*.
- Fu, Z.; Lam, W.; Yu, Q.; So, A. M.-C.; et al. 2023. Decoder-Only or Encoder-Decoder? Interpreting Language Model as a Regularized Encoder-Decoder. *arXiv preprint arXiv:2304.04052*.
- Gulati, A.; Qin, J.; Chiu, C.-C.; Parmar, N.; et al. 2020. Conformer: Convolution-augmented Transformer for speech recognition. In *Proc. INTERSPEECH*.
- Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *Proc. NeurIPS*.
- Holtzman, A.; Buys, J.; Du, L.; Forbes, M.; et al. 2019. The Curious Case of Neural Text Degeneration. In *Proc. ICLR*.
- Hsu, W.-N.; Zhang, Y.; Weiss, R. J.; Zen, H.; et al. 2019. Hierarchical generative modeling for controllable speech synthesis. In *Proc. ICLR*.
- Ippolito, D.; Kriz, R.; Sedoc, J.; Kustikova, M.; et al. 2019. Comparison of Diverse Decoding Methods from Conditional Language Models. In *Proc. ACL*.
- Jeong, M.; Kim, H.; Cheon, S. J.; Choi, B. J.; et al. 2021. Diff-TTS: A denoising diffusion model for text-to-speech. In *Proc. Interspeech*.
- Kahn, J.; Rivière, M.; Zheng, W.; Kharitonov, E.; et al. 2020. Libri-light: A benchmark for asr with limited or no supervision. In *Proc. ICASSP. IEEE*.
- Kharitonov, E.; Vincent, D.; Borsos, Z.; Marinier, R.; et al. 2023. Speak, Read and Prompt: High-fidelity text-to-speech with minimal supervision. *arXiv preprint arXiv:2302.03540*.
- Kim, J.; Kim, S.; Kong, J.; and Yoon, S. 2020. Glow-TTS: A generative flow for text-to-speech via monotonic alignment search. *Proc. NeurIPS*.
- Kim, J.; Lee, K.; Chung, S.; and Cho, J. 2024. CLaM-TTS: Improving Neural Codec Language Model for Zero-Shot Text-to-Speech. In *Proc. ICLR*.
- Le, M.; Vyas, A.; Shi, B.; Karrer, B.; et al. 2023. Voicebox: Text-Guided Multilingual Universal Speech Generation at Scale. In *Proc. NeurIPS*.
- Lee, Y.; Shin, J.; and Jung, K. 2022. Bidirectional variational inference for non-autoregressive text-to-speech. In *Proc. ICLR*.
- Lipman, Y.; Chen, R. T.; Ben-Hamu, H.; Nickel, M.; et al. 2023. Flow Matching for Generative Modeling. In *Proc. ICLR*.
- McAuliffe, M.; Socolof, M.; Mihuc, S.; Wagner, M.; et al. 2017. Montreal Forced Aligner: Trainable Text-Speech Alignment Using Kaldi. In *Proc. Interspeech*.
- Miao, C.; Liang, S.; Chen, M.; Ma, J.; et al. 2020. Flow-TTS: A non-autoregressive network for text to speech based on flow. In *Proc. ICASSP. IEEE*.
- Oord, A. v. d.; Dieleman, S.; Zen, H.; Simonyan, K.; et al. 2017. WaveNet: A generative model for raw audio. In *Proc. ICML*.
- Panayotov, V.; Chen, G.; Povey, D.; and Khudanpur, S. 2015. LibriSpeech: An ASR corpus based on public domain audio books. In *Proc. ICASSP. IEEE*.
- Ramesh, A.; Dhariwal, P.; Nichol, A.; Chu, C.; et al. 2022. Hierarchical Text-Conditional Image Generation with CLIP Latents. *arXiv preprint arXiv:2204.06125*.
- Rubenstein, P. K.; Asawaroengchai, C.; Nguyen, D. D.; Bapna, A.; et al. 2023. AudioPaLM: A Large Language Model That Can Speak and Listen. *arXiv preprint arXiv:2306.12925*.
- Shen, K.; Ju, Z.; Tan, X.; Liu, Y.; et al. 2023. NaturalSpeech 2: Latent diffusion models are natural and zero-shot speech and singing synthesizers. *arXiv preprint arXiv:2304.09116*.
- Sohl-Dickstein, J.; Weiss, E.; Maheswaranathan, N.; and Ganguli, S. 2015. Deep unsupervised learning using nonequilibrium thermodynamics. In *Proc. ICML. PMLR*.
- Song, Y.; and Ermon, S. 2020. Improved techniques for training score-based generative models. *Proc. NeurIPS*.
- Thoppilan, R.; De Freitas, D.; Hall, J.; Shazeer, N.; et al. 2022. LaMDA: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; et al. 2017a. Attention is all you need. *Proc. NeurIPS*.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; et al. 2017b. Attention is all you need. *Proc. NeurIPS*, 30.

Vyas, A.; Shi, B.; Le, M.; Tjandra, A.; Wu, Y.-C.; Guo, B.; Zhang, J.; Zhang, X.; Adkins, R.; Ngan, W.; et al. 2023. Audiobox: Unified audio generation with natural language prompts. *arXiv preprint arXiv:2312.15821*.

Wang, C.; Chen, S.; Wu, Y.; Zhang, Z.; et al. 2023a. Neural codec language models are zero-shot text to speech synthesizers. *arXiv preprint arXiv:2301.02111*.

Wang, J.; Du, Z.; Chen, Q.; Chu, Y.; et al. 2023b. LauraGPT: Listen, Attend, Understand, and Regenerate Audio with GPT. *arXiv preprint arXiv:2310.04673*.

Wang, X.; Thakker, M.; Chen, Z.; Kanda, N.; et al. 2023c. SpeechX: Neural codec language model as a versatile speech transformer. *arXiv preprint arXiv:2308.06873*.

Wang, Y.; Skerry-Ryan, R.; Stanton, D.; Wu, Y.; et al. 2017. Tacotron: Towards end-to-end speech synthesis. In *Proc. INTERSPEECH*.

Wu, S.; and Shi, Z. 2022. ItôWave: Itô stochastic differential equation is all you need for wave generation. In *Proc. ICASSP*. IEEE.

Yang, Z.; Dai, Z.; Yang, Y.; Carbonell, J.; et al. 2019. XLNet: Generalized autoregressive pretraining for language understanding. *Proc. NeurIPS*.

Zen, H.; Dang, V.; Clark, R.; Zhang, Y.; et al. 2019. LibriTTS: A Corpus Derived from LibriSpeech for Text-to-Speech. In *Proc. INTERSPEECH*.

Zen, H.; Tokuda, K.; and Black, A. W. 2009. Statistical parametric speech synthesis. *speech communication*, 51(11): 1039–1064.