

Edit Content, Preserve Acoustics: Imperceptible Text-Based Speech Editing via Self-Consistency Rewards

Yong Ren^{1,2}, Jiangyan Yi^{†3}, Jianhua Tao^{†3,4}, Zhengqi Wen^{†4}, Tao Wang¹

¹The State Key Laboratory of Multimodal Artificial Intelligence Systems,
Institute of Automation, Chinese Academy of Sciences, China

²School of Artificial Intelligence, University of Chinese Academy of Sciences, China

³Department of Automation, Tsinghua University, China

⁴Beijing National Research Center for Information Science and Technology, Tsinghua University, China
jhtao@tsinghua.edu.cn, yijy@tsinghua.edu.cn, zqwen@tsinghua.edu.cn

Abstract—Imperceptible text-based speech editing allows users to modify spoken content by altering the transcript. It demands that modified segments fuse seamlessly with the surrounding context. Prevalent methods operating in the acoustic space suffer from inherent content-style entanglement, leading to generation instability and boundary artifacts. In this paper, we propose a novel framework grounded in the principle of ‘Edit Content, Preserve Acoustics’. Our approach relies on two core components: (1) Structural Foundations, which decouples editing into a stable semantic space while delegating acoustic reconstruction to a Flow Matching decoder; and (2) Perceptual Alignment, which employs a novel Self-Consistency Rewards Group Relative Policy Optimization. By leveraging a pre-trained Text-to-Speech model as an implicit critic—complemented by strict intelligibility and duration constraints—we effectively align the edited semantic token sequence with the original context. Empirical evaluations demonstrate that our method significantly outperforms state-of-the-art autoregressive and non-autoregressive baselines, achieving superior intelligibility, robustness, and perceptual quality.

Index Terms—text-based speech editing, semantic token, reinforcement learning, self-consistency rewards

I. INTRODUCTION

Text-based speech editing enables the seamless modification of spoken content—inserting, deleting, or substituting words—by simply altering the transcript, eliminating the need for costly re-recording sessions [1]–[5]. This technology is invaluable for applications ranging from podcast correction and audiobook revision to post-production dialogue modification [6], [7].

Despite significant progress, achieving ‘imperceptible’ editing remains a formidable challenge. Early Non-Autoregressive (NAR) approaches [3], [8]–[10] offered inference stability but often failed to model long-range dependencies, resulting in monotonous prosody. Conversely, recent Autoregressive (AR) models based on Neural Codec Language Models (NCLMs) [6], [11]–[14] have achieved state-of-the-art (SOTA) naturalness. However, these methods typically operate on acoustic tokens [15], [16] where content and style are entangled. This coupling often compromises robustness, lead-

ing to hallucinations and boundary artifacts when modifying content [17].

To advance the state of text-based speech editing, we revisit the task and contend that it differs fundamentally from Text-to-Speech (TTS). We characterize text-based speech editing as a context-constrained incremental generation problem. Achieving imperceptibility requires balancing modification with continuity; to this end, we propose a framework grounded in the principle of *Edit Content, Preserve Acoustics*:

Structural Foundations for Acoustic Preservation. A primary limitation of prevalent text-based speech editing approaches lies in their direct manipulation of acoustic features. The tight coupling of linguistic content and acoustic timbre renders simultaneous prediction difficult and unstable, frequently resulting in hallucinations or boundary artifacts. To mitigate the interference of acoustic variations, we argue that text-based speech editing is best approached by modifying semantics first, followed by unified acoustic integration. Therefore, we decouple the editing process. We shift content manipulation to a disentangled semantic space, generating edited tokens that represent only linguistic content and coarse prosody. Subsequently, acoustic reconstruction is delegated to a Flow Matching [18] decoder. This hierarchical paradigm ensures that both the modified region and the original context are projected into a unified acoustic manifold, thereby preserving the fundamental acoustic coherence while allowing for precise content manipulation.

Perceptual Alignment for Further Coherence. Although structural decoupling maintains the acoustic foundation (such as timbre and acoustic environment), semantic tokens inherently encode rhythm and paralinguistic information. Consequently, ensuring that the edited content fuses indistinguishably with the utterance requires explicit perceptual alignment. Reinforcement Learning with Verifiable Rewards (RLVR) has emerged as a powerful paradigm for aligning Large Language Models (LLMs) [19]–[21], with initial explorations in speech generation [22], [23]. However, existing speech rewards typically target Text-to-Speech (TTS) metrics like intelligibility or

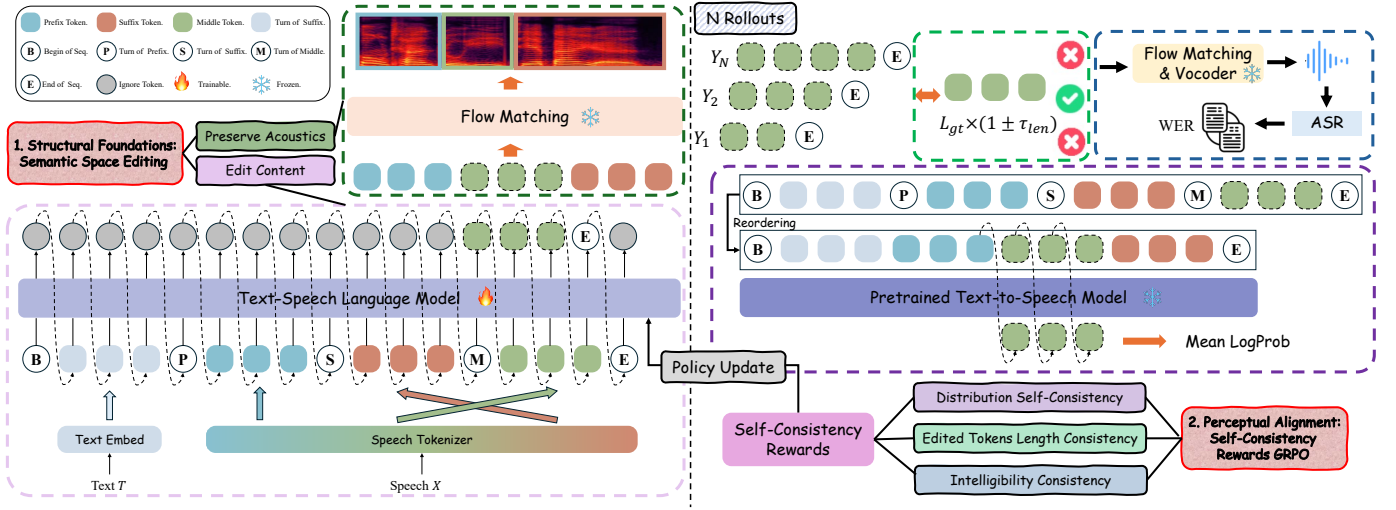


Fig. 1. The overall framework of our proposed method. The pipeline consists of two stages: (1) **Structural Foundations** (Left), where we employ a semantic-token-based LLM for conditional content infilling followed by a Flow Matching decoder for acoustic reconstruction; and (2) **Perceptual Alignment** (Right), where the policy is fine-tuned via Self-Consistency Rewards GRPO.

speaker similarity, failing to address the seamless fusion of the edited region with the surrounding context beyond mere timbre coherence. To address this, we introduce a novel alignment mechanism using a pre-trained Text-to-Speech (TTS) model as an implicit critic. Based on the premise that a powerful TTS model captures the distribution of natural speech, we utilize the conditional likelihood of the edited tokens given the context as a statistical proxy for the coherence between the edited region and the entire sentence. Complemented by strict constraints on Automatic Speech Recognition (ASR) Word Error Rate (WER) and duration validity, we construct a composite reward that aligns the generation with human perceptual expectations.

In summary, our contributions are as follows:

- We propose a novel framework for imperceptible text-based speech editing, addressing the problem through structural foundations and perceptual alignment.
- **Semantic Space Editing.** We adopt a semantic-token-based architecture to decouple content editing from acoustic reconstruction, significantly reducing artifacts compared to acoustic-token-based baselines.
- **Self-Consistency Rewards GRPO.** We are the first to utilize a pre-trained TTS model as a consistency critic within a Reinforcement Learning (RL) framework for speech editing. Combined with ASR and length rewards, it effectively enhances global coherence and achieves perceptual alignment.
- Empirical evaluations demonstrate that our method significantly outperforms NAR and AR baselines, achieving superior intelligibility, robustness and perceptual quality.

II. METHOD

This section details the proposed framework for imperceptible text-based speech editing. We first outline the overall architecture, followed by the specific implementation of **Structural Foundations** and the **Perceptual Alignment** mechanism.

A. Overall architecture

Guided by the core principle of *Edit Content, Preserve Acoustics*, we design a hierarchical framework as illustrated in Figure 1. Our approach orchestrates this principle through two stages:

- **Structural Foundations via Decoupling:** To insulate acoustic stability from content modification, we restrict editing operations to a discrete semantic space. This effectively isolates linguistic manipulation from acoustic timbre variations, which are subsequently reconstructed by a Flow Matching decoder within a unified manifold.
- **Perceptual Alignment via RL:** To ensure the edited region fuses indistinguishably with the bidirectional context, we introduce a *Self-Consistency Rewards GRPO* stage. We leverage a pre-trained TTS model as an implicit critic to maximize the conditional likelihood of the generated tokens, thereby enhancing global coherence of semantic tokens without requiring paired ground-truth data.

B. Structural Foundations: Semantic Space Editing

The primary impediment to structural preservation in prior AR approaches lies in the inherent entanglement within acoustic codecs [6]. Modifying tokens in this entangled acoustic space inevitably perturbs the style trajectory, causing boundary artifacts. To address this, we adopt a decoupled architecture comprising an LLM for semantic generation and a Flow Matching decoder for acoustic reconstruction [24].

We formulate the text-based speech editing task as a conditional token infilling problem within the semantic space. To facilitate this, we employ a Prefix-Suffix-Middle (PSM) formatting strategy. Given the raw audio X , we first utilize a semantic tokenizer to encode it into a discrete token sequence S . We then partition S into three segments— S_{pre} , S_{mid} , and S_{suf} . As illustrated in Figure 1 (left), to construct the input

sequence \mathbf{Q} , we concatenate the encoded target text with the semantic tokens \mathbf{S}_{pre} and \mathbf{S}_{suf} , separating each modality and region with special tokens:

$$\mathbf{Q} = [\text{Enc}(\mathbf{T}); \mathbf{S}_{\text{pre}}; \mathbf{S}_{\text{suf}}], \quad \mathbf{V} = \mathbf{S}_{\text{mid}}, \quad (1)$$

where \mathbf{V} is the target to predict and \mathbf{S}_{mid} represents the semantic tokens corresponding to the edited region.

We employ a decoder-only transformer as the policy model π_θ . During this supervised training stage, we optimize the standard negative log-likelihood objective focusing on the prediction of the missing middle tokens \mathbf{S}_{mid} , conditioned on the bidirectional context \mathbf{Q} :

$$\mathcal{L}(\theta) = - \sum_{t=1}^{|\mathbf{S}_{\text{mid}}|} \log \pi_\theta(s_{\text{mid},t} \mid \mathbf{Q}, s_{\text{mid},<t}). \quad (2)$$

This formulation ensures structural integrity: the model learns to ‘infill’ the semantic content that bridges the prefix and suffix seamlessly, while the subsequent Flow Matching decoder and vocoder reconstructs the entire waveform in a unified acoustic manifold.

C. Perceptual Alignment: Self-Consistency Rewards GRPO

While editing in the semantic space ensures fundamental acoustic coherence, semantic tokens encapsulate prosodic and paralinguistic information. Consequently, autoregressive sampling can still yield hallucinations or prosodic mismatches, resulting in perceptual discontinuity. To achieve imperceptibility—where the edited region fuses indistinguishably with the unedited context—we propose perceptual alignment. This mechanism ensures that the generated tokens statistically align with the surrounding context under the distribution of natural speech. To this end, we introduce a novel *Self-Consistency Rewards GRPO*, as illustrated in Figure 1 (right).

Unlike traditional actor-critic frameworks that require a computationally expensive value network, GRPO estimates the baseline directly from the collective performance of multiple outputs sampled for the same prompt. For each editing query \mathbf{Q} , we sample a group of G candidate sequences $\{\mathbf{o}_1, \dots, \mathbf{o}_G\}$ from the old policy $\pi_{\theta_{\text{old}}}$. The core of this optimization lies in the computation of the Relative Advantage \hat{A}_i , which reflects the merit of an individual sample relative to its peers. Formally, for each output \mathbf{o}_i with its corresponding reward R_i , the advantage is computed as:

$$\hat{A}_i = \frac{R_i - \text{mean}(R_j)}{\text{std}(R_j)}, \quad (3)$$

where i indexes the i -th completion among G candidates for the same input \mathbf{Q} .

To bridge the gap between discrete semantic editing and continuous acoustic realization, we design specific rewards that enforce statistical consistency with natural speech.

Log-Probability Self-Consistency Reward (r_{sc}): This component serves as the cornerstone of our perceptual alignment strategy. We leverage the pre-trained TTS model itself as an *implicit critic* of naturalness.

Formulation: For a generated token sequence $\hat{\mathbf{S}}_{\text{mid}}$ by speech editing model, we compute the average log-likelihood under the frozen TTS reference model π_{tts} :

$$r_{\text{sc}} = \frac{1}{|\hat{\mathbf{S}}_{\text{mid}}|} \sum_{t=1}^{|\hat{\mathbf{S}}_{\text{mid}}|} \log \pi_{\text{tts}}(\hat{s}_{\text{mid},t} \mid [\text{Enc}(\mathbf{T}); \mathbf{S}_{\text{pre}}], \hat{s}_{\text{mid},<t}). \quad (4)$$

Theoretical Justification: The pre-trained TTS model π_{tts} approximates the true distribution of natural speech P_{speech} . Maximizing the expected reward $J(\theta) = \mathbb{E}_{\hat{s} \sim \pi_\theta} [r_{\text{sc}}(\hat{s})]$ is equivalent to minimizing the Cross-Entropy between the policy’s generation and the TTS model’s prior:

$$\begin{aligned} \max_{\theta} J(\theta) &\iff \\ \min_{\theta} H(\pi_\theta, \pi_{\text{tts}}) &= \min_{\theta} [H(\pi_\theta) + \mathbb{D}_{\text{KL}}(\pi_\theta \parallel \pi_{\text{tts}})], \end{aligned} \quad (5)$$

where $H(\pi_\theta, \pi_{\text{tts}}) = -\mathbb{E}_{x \sim \pi_\theta} [\log \pi_{\text{tts}}(x)]$. By maximizing r_{sc} , we strictly constrain the policy π_θ to stay within the *high-probability manifold* of natural speech defined by π_{tts} . This ensures that the edited region maintains the same prosodic coherence, rhythm, and co-articulation patterns as the surrounding unedited context, effectively preventing the out-of-distribution artifacts common in editing.

Intelligibility Reward (r_{wer}): Relying exclusively on r_{sc} renders the policy susceptible to reward hacking. Since high-probability tokens often correspond to common patterns (e.g., silence or simple repetitions), the model may converge to degenerate local optima that maximize likelihood but lack semantic content. To mitigate such mode collapse and strictly enforce content accuracy, we incorporate an intelligibility penalty:

$$r_{\text{wer}} = 1 - \text{WER}(\mathbf{W}_{\text{rec}}, \mathbf{T}_{\text{tgt}}), \quad (6)$$

where WER is computed by an ASR model on the waveform \mathbf{W}_{rec} reconstructed from the tokens. This explicitly penalizes repetitive or unintelligible outputs, forcing the model to balance naturalness with linguistic precision.

Gated Self-Consistency Rewards Fusion To orchestrate the trade-off between acoustic naturalness (r_{sc}) and linguistic accuracy (r_{wer}), we propose a *Gated Reward Aggregation* strategy. A simple linear combination of rewards risks allowing the policy to exploit one metric at the expense of the other. To mitigate this, we impose a hard validity constraint. We define the total reward R as a piecewise function that zeros out the feedback for samples failing to meet minimal quality standards:

$$R = \begin{cases} R_{\text{base}} + r_{\text{sc}} + r_{\text{wer}} & \text{if } \mathbb{I}_{\text{valid}} \\ 0 & \text{otherwise} \end{cases}, \quad (7)$$

where R_{base} is a shaping constant to ensure positive rewards for valid samples.

The validity indicator $\mathbb{I}_{\text{valid}}$ acts as a binary filter, ensuring the generated speech satisfies both intelligibility and duration consistency constraints:

$$\mathbb{I}_{\text{valid}} = \underbrace{(\text{WER}(\mathbf{W}_{\text{rec}}, \mathbf{T}_{\text{tgt}}) \leq \tau_{\text{wer}})}_{\text{Content Integrity}} \wedge \underbrace{\left(\left| \frac{L_{\text{gen}} - L_{\text{tgt}}}{L_{\text{tgt}}} \right| \leq \tau_{\text{len}} \right)}_{\text{Duration Stability}}, \quad (8)$$

TABLE I
PERFORMANCE COMPARISON ON THE TEXT-BASED SPEECH EDITING BENCHMARK. THE SYMBOL '◊' INDICATES THAT THE RESULTS ARE CITED DIRECTLY FROM THE ORIGINAL PAPER [14]. **RED** INDICATES THE BEST RESULT, AND **BLUE** INDICATES THE SECOND BEST.

| Edit Type | Model | Performance | | | | | | | |
|--------------|----------------|----------------------|-------------|-------------------|-------------|----------------------|-------------|-------------------|-------------|
| | | WER(%)↓ basic full | | SIM↑ basic full | | DNSMOS↑ basic full | | MOS↑ basic full | |
| Insertion | FluentSpeech | 12.00 | 11.91 | 0.60 | 0.60 | 2.90 | 2.91 | 3.42 | 3.41 |
| | VoiceCraft | 10.70 | 12.94 | 0.67 | 0.67 | 3.00 | 3.00 | 3.65 | 3.64 |
| | Ming-UniAudio◊ | 6.63 | 7.59 | 0.79 | 0.79 | - | - | - | - |
| | Ours | 4.70 | 5.12 | 0.82 | 0.82 | 3.14 | 3.13 | 3.92 | 3.90 |
| | Ours (w. GRPO) | 4.50 | 4.97 | 0.82 | 0.82 | 3.17 | 3.18 | 4.08 | 4.06 |
| Deletion | FluentSpeech | 8.16 | 8.78 | 0.51 | 0.52 | 2.91 | 2.91 | 3.45 | 3.44 |
| | VoiceCraft | 16.99 | 17.88 | 0.60 | 0.62 | 3.01 | 3.04 | 3.28 | 3.30 |
| | Ming-UniAudio◊ | 14.85 | 27.60 | 0.76 | 0.74 | - | - | - | - |
| | Ours | 7.38 | 7.70 | 0.76 | 0.77 | 3.07 | 3.08 | 3.85 | 3.86 |
| | Ours (w. GRPO) | 6.91 | 6.88 | 0.77 | 0.78 | 3.09 | 3.09 | 3.96 | 3.95 |
| Substitution | FluentSpeech | 4.66 | 4.65 | 0.51 | 0.51 | 2.92 | 2.92 | 3.48 | 3.47 |
| | VoiceCraft | 11.98 | 12.73 | 0.58 | 0.59 | 3.01 | 3.02 | 3.60 | 3.61 |
| | Ming-UniAudio◊ | 8.99 | 7.64 | 0.78 | 0.77 | - | - | - | - |
| | Ours | 4.40 | 4.61 | 0.78 | 0.78 | 3.12 | 3.09 | 3.94 | 3.92 |
| | Ours (w. GRPO) | 4.13 | 4.41 | 0.78 | 0.78 | 3.15 | 3.11 | 4.05 | 4.03 |

where $L_{\text{gen}} = |\hat{\mathbf{S}}_{\text{mid}}|$ and $L_{\text{gt}} = |\mathbf{S}_{\text{mid}}|$. The thresholds τ_{wer} and τ_{len} define the tolerance for admissible errors. This gating mechanism effectively prunes the optimization space, preventing the policy from learning from catastrophic failures, thereby stabilizing the GRPO training process.

III. EXPERIMENTS

A. Experimental Settings

Datasets. We conduct training on the Libriheavy [25] dataset, a large-scale corpus comprising approximately 50,000 hours of English speech from LibriVox. For evaluation, we establish two distinct benchmarks to rigorously assess editing performance. The first, adapted from the Ming-Freeform-Audio-Edit-Benchmark [14] (basic & full), targets precise editing operations (Insertion, Deletion, Substitution). To ensure precision, we define ground-truth editing intervals via forced alignment. Specifically, WhisperX [26] is employed to align original and target texts with the audio, automatically extracting accurate timestamps for editing masks. To evaluate stability under varying edit lengths, we constructed a supplementary test set derived from Seed-TTS-Eval [27]. We randomly sampled 200 utterances (>4s) and applied random editing masks with fixed durations of $\{0.5\text{s}, \dots, 2.5\text{s}\}$.

Baselines. We benchmark our framework against three representative models covering diverse editing paradigms:

- **FluentSpeech** [28]: A NAR diffusion-based model generating mel-spectrograms directly.
- **VoiceCraft** [6]: A SOTA AR NCLM operating on quantized acoustic tokens.
- **Ming-UniAudio** [14]: A unified LLM designed for speech understanding, editing, and generation.

Evaluation Metrics. We employ both objective and subjective metrics for a comprehensive evaluation: (1) **WER**: Calculated

on the full edited utterance using Whisper [29] to assess global intelligibility and boundary fusion quality. (2) **Speaker Similarity (SIM)**: The cosine similarity between WavLM [30] embeddings of the edited and original speech, measuring timbre preservation. (3) **DNSMOS** [31]: A neural perceptual metric estimating speech naturalness and quality. (4) **Mean Opinion Score (MOS)**: To assess perceptual imperceptibility, we conducted subjective listening tests. We randomly selected 90 samples and recruited 8 raters to evaluate the naturalness and coherence on a scale of 1 (Bad) to 5 (Excellent).

B. Implementation Details

Our pipeline comprises two stages: Supervised Pre-training and RL Alignment. We adopt the semantic tokenizer, Flow Matching decoder, and HiFiGAN [32] vocoder from CosyVoice3 [24], keeping these components frozen. We trained the Semantic Token LLM on Libriheavy. The training configuration included a learning rate of 1×10^{-5} , gradient accumulation steps of 2, and a maximum of 10 epochs. During RL training, we set the learning rate to 1×10^{-6} , batch size to 4, and group size (rollout) to 8. The KL divergence coefficient β was set to 0.01. Training proceeded for 400 steps with gradient accumulation steps of 10. The Log-Probability Reward was computed by the pre-trained TTS model CosyVoice3 [24], and the ASR Reward by SenseVoiceSmall [33]. τ_{wer} and τ_{len} are set to 0.2. Our experiments were conducted on 8 NVIDIA H800 GPUs.

C. Main Results Analysis

Table I presents the performance comparison on the first benchmark. Audio samples are available here.¹

¹<https://icme26-speech-editing.netlify.app>.

TABLE II
ROBUSTNESS EVALUATION ON THE SUBSET OF SEED-TTS TEST SET
ACROSS VARYING MASKED DURATIONS (0.5s TO 2.5s). **RED** INDICATES
THE BEST RESULT, AND **BLUE** INDICATES THE SECOND BEST.

| Metrics | Model | Masked Duration | | | | |
|----------|----------------|-----------------|--------------|--------------|--------------|--------------|
| | | 0.5s | 1s | 1.5s | 2s | 2.5s |
| WER(%) ↓ | FluentSpeech | 7.533 | 7.107 | 7.231 | 8.300 | 7.390 |
| | VoiceCraft | 8.504 | 10.505 | 10.813 | 11.525 | 11.190 |
| | Ours | 3.342 | 3.858 | 4.126 | 4.430 | 4.333 |
| | Ours (w. GRPO) | 3.202 | 3.400 | 4.067 | 4.117 | 4.227 |
| SIM ↑ | FluentSpeech | 0.797 | 0.750 | 0.685 | 0.615 | 0.535 |
| | VoiceCraft | 0.790 | 0.761 | 0.723 | 0.695 | 0.639 |
| | Ours | 0.866 | 0.855 | 0.840 | 0.828 | 0.809 |
| | Ours (w. GRPO) | 0.865 | 0.854 | 0.840 | 0.829 | 0.811 |
| DNSMOS ↑ | FluentSpeech | 2.915 | 2.930 | 2.975 | 2.991 | 3.006 |
| | VoiceCraft | 2.970 | 2.978 | 3.016 | 3.021 | 3.008 |
| | Ours | 3.126 | 3.139 | 3.138 | 3.138 | 3.128 |
| | Ours (w. GRPO) | 3.124 | 3.138 | 3.143 | 3.148 | 3.148 |

Intelligibility and Robustness. Our method, based on semantic token editing, achieves the lowest WER across all three editing operations, outperforming SOTA baselines. This structural advantage stems from decoupling semantic generation from acoustic rendering, which simplifies the modeling task. Furthermore, applying GRPO yields further WER reductions across all tasks. We attribute this to the ASR reward and the length constraint within the GRPO framework, which effectively penalize unintelligible content and gibberish.

Detailed analysis of different types of editing operations:

- **Insertion:** AR methods generally outperform NAR methods, with Ming-UniAudio also demonstrating competitive performance. The NAR baseline (FluentSpeech) exhibits the highest WER. As NAR models rely on mask prediction, the duration of the masked region often mismatches the length required for the new text in insertion tasks, leading to severe artifacts.
- **Deletion:** This task proves most challenging for traditional AR models (e.g., VoiceCraft), which suffer from severe hallucinations and often fail to emit the EOS token in time. Conversely, NAR models perform better by simply predicting silence. Our method achieves the best performance; by editing only semantic tokens, we significantly reduce prediction difficulty, enabling precise stopping. Moreover, with GRPO, the WER drops drastically (0.47% on basic and 0.82% on full), confirming that the length reward effectively suppresses repetition and incoherent generation.
- **Substitution:** While NAR methods excel here due to similar durations between original and edited text, our method still surpasses the strong NAR baseline.

Speaker Similarity and Perceptual Quality. Our method consistently achieves the highest SIM, DNSMOS, and subjective MOS scores, outperforming all baselines.

In terms of *SIM*, the Audio-LLM-based Ming-UniAudio scores slightly lower than ours, followed by VoiceCraft, with FluentSpeech performing worst. This validates the superiority

of AR-based approaches over diffusion-based NAR methods in capturing speaker characteristics. Notably, GRPO does not significantly alter the *SIM* score. This is expected, as we optimize the policy for semantic tokens, while timbre preservation is primarily handled by the fixed Flow Matching decoder.

Regarding *DNSMOS* and Subjective *MOS*, our method outperforms baselines in naturalness.

- Our method significantly outperforms both the acoustic AR baseline (VoiceCraft) and the NAR baseline (FluentSpeech).
- The introduction of GRPO yields substantial improvements in perceptual metrics. This is driven by the Self-Consistency Rewards, which ensure the prosody of the edited region aligns seamlessly with the unedited context, rendering the edit ‘imperceptible’.
- Subjective *MOS* results mirror the objective *DNSMOS*, confirming that human listeners perceive our method—especially with GRPO alignment—as the most natural.

D. Robustness on Edited Duration

Table II illustrates the performance stability as the masked duration increases from 0.5s to 2.5s.

Impact on Intelligibility. Our method consistently achieves the lowest WER across all durations, with GRPO providing further reductions. Unlike the main benchmark, in this setup (where mask duration equals generated text duration), the NAR baseline (FluentSpeech) outperforms the traditional AR baseline (VoiceCraft). As the edited duration increases, VoiceCraft’s WER escalates sharply, reflecting the error accumulation inherent in acoustic token autoregression. In contrast, our method exhibits a much slower rate of degradation. Even at 2.5s, our WER remains significantly lower than FluentSpeech, demonstrating robust long-context modeling.

Impact on Speaker Similarity. Our method maintains the highest speaker similarity across all durations. Consistent with the main results, GRPO has a negligible effect on *SIM*. Among baselines, VoiceCraft outperforms FluentSpeech. Crucially, the NAR method shows a drastic decline in similarity as duration increases (dropping to 0.535 at 2.5s). Traditional AR also degrades, albeit more slowly. Our method, however, maintains high similarity with minimal decay. This robustness results directly from our decoupled architecture: the Flow Matching decoder performs a unified acoustic reconstruction, rendering it insensitive to the length of the edited semantic tokens.

Impact on Naturalness (DNSMOS). Our method consistently achieves the best naturalness scores, with the impact of GRPO becoming increasingly pronounced as edit length grows. For short edits (e.g., 0.5s), the influence of GRPO is marginal. However, as the duration extends to 2.5s, the gap widens significantly. While baseline methods and our non-aligned model show little improvement or plateau, Ours (with GRPO) shows a clear upward trend. This validates the efficacy of our *Self-Consistency Reward*: by leveraging a pre-trained TTS model to approximate the natural speech distribution, RL

optimization guides the model to generate coherent prosody even for complex, long-form edits.

IV. CONCLUSION

In this paper, we presented a novel framework for imperceptible text-based speech editing through the principle of *Edit Content, Preserve Acoustics*. By shifting the editing operation from the acoustic space to a disentangled semantic space, we established a robust structural foundation for acoustic preservation. Furthermore, to ensure the edited region fuses indistinguishably with the context, we introduced a Perceptual Alignment stage via Self-Consistency Rewards GRPO. Extensive evaluations on two benchmarks demonstrate that our method significantly outperforms SOTA AR and NAR baselines, achieving superior intelligibility, speaker similarity, and perceptual naturalness, even in long-duration scenarios. Future work will explore extending this semantic-based framework and the GRPO alignment method to freeform speech editing.

REFERENCES

- [1] Zeyu Jin, Gautham J Mysore, Stephen Diverdi, Jingwan Lu, and Adam Finkelstein, “Voco: Text-based insertion and replacement in audio narration,” *ACM Transactions on Graphics (TOG)*, vol. 36, no. 4, pp. 1–13, 2017.
- [2] Daxin Tan, Liqun Deng, Yu Ting Yeung, Xin Jiang, Xiao Chen, and Tan Lee, “Editspeech: A text based speech editing system using partial inference and bidirectional fusion,” in *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2021, pp. 626–633.
- [3] Tao Wang, Jiangyan Yi, Ruibo Fu, Jianhua Tao, and Zhengqi Wen, “Campnet: Context-aware mask prediction for end-to-end text-based speech editing,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 2241–2254, 2022.
- [4] Tao Wang, Jiangyan Yi, Liqun Deng, Ruibo Fu, Jianhua Tao, and Zhengqi Wen, “Context-aware mask prediction network for end-to-end text-based speech editing,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 6082–6086.
- [5] Tao Wang, Jiangyan Yi, Ruibo Fu, Jianhua Tao, Zhengqi Wen, and Chu Yuan Zhang, “Emotion selectable end-to-end text-based speech editing,” *Artificial Intelligence*, vol. 329, pp. 104076, 2024.
- [6] Puyuan Peng, Po-Yao Huang, Shang-Wen Li, Abdelrahman Mohamed, and David Harwath, “Voicecraft: Zero-shot speech editing and text-to-speech in the wild,” in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2024, pp. 12442–12462.
- [7] Taewoo Kim, Uijong Lee, Hayoung Park, Choongsang Cho, Nam In Park, and Young Han Lee, “Instance-specific test-time training for speech editing in the wild,” *arXiv preprint arXiv:2506.13295*, 2025.
- [8] He Bai, Renjie Zheng, Junkun Chen, Mingbo Ma, Xintong Li, and Liang Huang, “A³t: Alignment-aware acoustic and text pretraining for speech synthesis and editing,” in *International Conference on Machine Learning*. PMLR, 2022, pp. 1399–1411.
- [9] Tao Wang, Jiangyan Yi, Ruibo Fu, Chunyu Qiang, Dading Chong, Chao Wang, Z Wen, J Tao, et al., “Speechpalette: A comprehensive speech editing method for text-based speech editing, one-shot tts and attributes editing,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025.
- [10] Rui Liu, Jiatian Xi, Ziyue Jiang, and Haizhou Li, “Fluenteditor2: Text-based speech editing by modeling multi-scale acoustic and prosody consistency,” *IEEE Transactions on Audio, Speech and Language Processing*, 2025.
- [11] Zhisheng Zheng, Puyuan Peng, Anuj Diwan, Cong Phuoc Huynh, Xiaohang Sun, Zhu Liu, Vimal Bhat, and David Harwath, “Voicecraft-x: Unifying multilingual, voice-cloning speech synthesis and speech editing,” in *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, 2025, pp. 2737–2756.
- [12] Baher Mohammad, Magaiya Zhussip, and Stamatios Lefkimmiatis, “Speak, edit, repeat: High-fidelity voice editing and zero-shot tts with cross-attentive mamba,” *arXiv preprint arXiv:2510.04738*, 2025.
- [13] Dongchao Yang, Jinchuan Tian, Xu Tan, Rongjie Huang, Songxiang Liu, Xuankai Chang, Jiatong Shi, Sheng Zhao, Jiang Bian, Xixin Wu, et al., “Uniaudio: An audio foundation model toward universal audio generation,” *arXiv preprint arXiv:2310.00704*, 2023.
- [14] Canxiang Yan, Chunxiang Jin, Dawei Huang, Haibing Yu, Han Peng, Hui Zhan, Jie Gao, Jing Peng, Jingdong Chen, Jun Zhou, et al., “Ming-uniaudio: Speech llm for joint understanding, generation and editing with unified representation,” *arXiv preprint arXiv:2511.05516*, 2025.
- [15] Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi, “High fidelity neural audio compression,” *Trans. Mach. Learn. Res.*, 2023.
- [16] Rithesh Kumar, Prem Seetharaman, Alejandro Luebs, Ishaan Kumar, and Kundan Kumar, “High-fidelity audio compression with improved rvqgan,” *Advances in Neural Information Processing Systems*, vol. 36, pp. 27980–27993, 2023.
- [17] Sung-Feng Huang, Heng-Cheng Kuo, Zhehuai Chen, Xuesong Yang, Pin-Jui Ku, Ante Jukic, Chao-Han Huck Yang, Yu Tsao, Yu-Chiang Frank Wang, Hung-yi Lee, et al., “Voicenong: Robust high-quality speech editing model without hallucinations,” in *Proc. Interspeech 2025*, 2025, pp. 3469–3473.

- [18] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le, “Flow matching for generative modeling,” in The Eleventh International Conference on Learning Representations, 2023.
- [19] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al., “Deepseekmath: Pushing the limits of mathematical reasoning in open language models,” arXiv preprint arXiv:2402.03300, 2024.
- [20] Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins, “Solving math word problems with process-and outcome-based feedback,” arXiv preprint arXiv:2211.14275, 2022.
- [21] Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe, “Let’s verify step by step,” in The Twelfth International Conference on Learning Representations, 2023.
- [22] Dong Zhang, Zhaowei Li, Shimin Li, Xin Zhang, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu, “Speechalign: Aligning speech generation to human preferences,” Advances in Neural Information Processing Systems, vol. 37, pp. 50343–50360, 2024.
- [23] Yicheng Zhong, Peiji Yang, and Zhisheng Wang, “Multi-reward grpo for stable and prosodic single-codebook tts llms at scale,” arXiv preprint arXiv:2511.21270, 2025.
- [24] Zhihao Du, Changfeng Gao, Yuxuan Wang, Fan Yu, Tianyu Zhao, Hao Wang, Xiang Lv, Hui Wang, Chongjia Ni, Xian Shi, et al., “Cosyvoice 3: Towards in-the-wild speech generation via scaling-up and post-training,” arXiv preprint arXiv:2505.17589, 2025.
- [25] Wei Kang, Xiaoyu Yang, Zengwei Yao, Fangjun Kuang, Yifan Yang, Liyong Guo, Long Lin, and Daniel Povey, “Libriheavy: A 50,000 hours asr corpus with punctuation casing and context,” in ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2024, pp. 10991–10995.
- [26] Max Bain, Jaesung Huh, Tengda Han, and Andrew Zisserman, “Whisperx: Time-accurate speech transcription of long-form audio,” in Proc. Interspeech 2023, 2023, pp. 4489–4493.
- [27] Philip Anastassiou, Jiawei Chen, Jitong Chen, Yuanzhe Chen, Zhuo Chen, Ziyi Chen, Jian Cong, Lelai Deng, Chuang Ding, Lu Gao, et al., “Seed-tts: A family of high-quality versatile speech generation models,” arXiv preprint arXiv:2406.02430, 2024.
- [28] Ziyue Jiang, Qian Yang, Jialong Zuo, Zhenhui Ye, Rongjie Huang, Yi Ren, and Zhou Zhao, “Fluentspeech: Stutter-oriented automatic speech editing with context-aware diffusion models,” in Findings of the Association for Computational Linguistics: ACL 2023, 2023, pp. 11655–11671.
- [29] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever, “Robust speech recognition via large-scale weak supervision,” in International conference on machine learning. PMLR, 2023, pp. 28492–28518.
- [30] Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al., “Wavlm: Large-scale self-supervised pre-training for full stack speech processing,” IEEE Journal of Selected Topics in Signal Processing, vol. 16, no. 6, pp. 1505–1518, 2022.
- [31] Chandan KA Reddy, Vishak Gopal, and Ross Cutler, “Dnsmos: A non-intrusive perceptual objective speech quality metric to evaluate noise suppressors,” in ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2021, pp. 6493–6497.
- [32] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae, “Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis,” Advances in neural information processing systems, vol. 33, pp. 17022–17033, 2020.
- [33] Keyu An, Qian Chen, Chong Deng, Zhihao Du, Changfeng Gao, Zhifu Gao, Yue Gu, Ting He, Hangrui Hu, Kai Hu, et al., “Funaudiollm: Voice understanding and generation foundation models for natural interaction between humans and llms,” arXiv preprint arXiv:2407.04051, 2024.