

WhispSynth: Scaling Multilingual Whisper Corpus through Real Data Curation and A Novel Pitch-free Generative Framework

Anonymous ACL submission

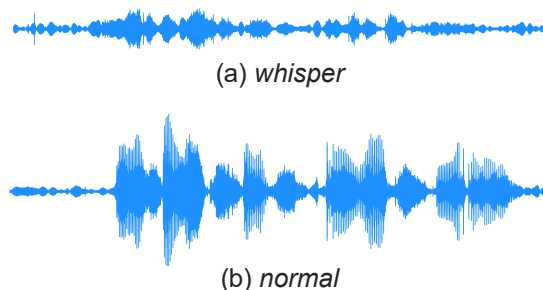
Abstract

Whisper generation is constrained by the difficulty of data collection. Because whispered speech has low acoustic amplitude, high-fidelity recording is challenging. In this paper, we introduce **WhispSynth**, a large-scale multilingual corpus constructed via a novel high-fidelity generative framework. Specifically, we propose a pipeline integrating Differentiable Digital Signal Processing (DDSP)-based pitch-free method with Text-to-Speech (TTS) models. This framework refines a comprehensive collection of resources, including our newly constructed WhispXXX dataset¹, into 118 hours of high-fidelity whispered speech from 479 speakers. Unlike standard synthetic or noisy real data, our data engine faithfully preserves source vocal timbre and linguistic content while ensuring acoustic consistency, providing a robust foundation for text-to-whisper research. Experimental results demonstrate that WhispSynth exhibits significantly higher quality than existing corpora. Moreover, our CosyWhisper, tuned with WhispSynth, achieves speech naturalness on par with ground-truth samples. We will release the implementation code and available resources to enable the reproduction of the WhispSynth generation pipeline.

1 Introduction

Whisper represents an innate and ubiquitous mode of human vocalisation. Primarily employed to attenuate acoustic propagation for secure communication (Tartter, 1989; Andersen, 2015; Rekimoto, 2023; Hiraki and Rekimoto, 2025), this modality also carries significant paralinguistic weight, often signalling intimacy or inducing a soothing effect. Notably, these affective characteristics have driven the recent surge in Autonomous Sensory Meridian Response (ASMR) content across social media.

¹We refer to the dataset as "WhispXXX" to maintain institutional anonymity. The real dataset name and resource will be disclosed upon publication.



tr: "When the sunlight strikes raindrops in the air, ..."

Figure 1: Different Dynamic Range. Whispers exhibit a significantly lower sound pressure level compared to normal speech, even when the linguistic content is identical.

Research interest in whisper generation is rapidly growing, yet a critical data bottleneck fundamentally constrains progress. Publicly available whispered corpora are typically characterized by limited scale, heterogeneous recording quality, and a lack of speaker diversity. Specifically, the inherently low sound pressure level of whispers (see Figure 1) makes high-fidelity capture exceptionally difficult. Furthermore, existing datasets often contain substantial acoustic artifacts that are resistant to standard speech enhancement, because whispered speech is produced with predominantly noise-driven excitation.

While Text-to-Speech (TTS) remains the primary paradigm for synthesising whisper, and general TTS has made remarkable strides in capturing nuanced prosodic features such as rhythm, intonation, and even non-verbal vocalizations (Du et al., 2025; Zhou et al., 2025; Chen et al., 2025), achieving high-fidelity whisper generation remains a significant challenge. Recently, we have observed that commercial services such as Doubao AI Chatbot (ByteDance, 2025), ElevenLabs TTS (ElevenLabs, 2025), and MiniMax TTS (ElevenLabs,

2025) can produce natural-sounding whispered speech. However, these systems rely on proprietary datasets and closed-source pipelines. Consequently, the lack of high-quality, open-access resources remains a critical bottleneck that directly impedes advances in general whisper synthesis.

To further investigate the limitations of current open-source TTS models in whisper generation, we reveal that TTS models remain severely flawed at zero-shot whisper generation. Even when provided with explicit whisper prompts or textual instructions, these models fail to bridge the domain gap, highlighting a significant deficiency in serving this special modality. This fragility can be attributed to two complementary biases: (i) corpora and vocoders are predominantly optimized for modal phonation, establishing implicit priors that favor harmonic-rich signals with clear fundamental frequencies (F0); and (ii) speaker-embedding extractors trained on such data frequently misconstrue whisper prompts as hoarseness or environmental noise rather than authentic breathy voice.

Paradoxically, these TTS models trained on normal speech possess a latent capacity for whisper generation, derived from sporadic devoiced segments that share similar acoustic energy distributions. The overlap is anatomically grounded: whispered production suppresses vocal-fold vibration while retaining turbulent airflow at the same supralaryngeal constriction sites exploited for voiceless obstruent in normal speech (Ito et al., 2005). Consequently, the acoustic realisations of whisper and devoiced segments are governed by nearly identical aerodynamic-articulatory constraints, endowing standard models with a dormant yet exploitable reservoir of whisper-like priors.

To address these limitations, we introduce **WhispSynth**, a large-scale, multilingual corpus constructed via a novel high-fidelity generative framework. Specifically, we propose a whisper generation pipeline integrating Differentiable Digital Signal Processing (DDSP)-based pitch-free with TTS models. This framework refines a comprehensive collection of resources into 118 hours of high-fidelity 24 kHz whispers over 479 speakers. The main contributions are summarized as follows:

1. **WhispReal**: We construct the first union of realistic whisper by consolidating 6 publicly available corpora and releasing our own newly recorded Mandarin dataset, **WhispXXX** (~45 h). The resulting collection, termed **Whis-**

pReal (~118 h), provides standardized splits and metadata, establishing a robust data foundation for the field.

2. **WhispSynth**: We propose a scalable data engine that synergizes state-of-the-art TTS (e.g., CosyVoice 3) with a novel pitch-free method. This pipeline effectively disentangles pitch from wav signal, transforming the noisy **WhispReal** source into **WhispSynth**—a studio-grade synthetic corpus of approximately 118 hours. The generated whisper preserves linguistic content, reducing CER and WER by 11%, and also improves audio quality, increasing DNSMOS by 3%.
3. **CosyWhisper**: We fine-tune CosyVoice 3 on our **WhispSynth** to develop **CosyWhisper**. Experimental results demonstrate that training on **WhispSynth** yields significantly better performance than training on realistic recordings, achieving speech naturalness on par with ground-truth samples. We will release the fine-tuned model as the first open-source and reproducible *text-to-whisper* system.

2 WhispReal

WhispReal dataset is a curated collection of existing publicly available whispered speech corpora combined with our proprietary **WhispXXX (anonymized)** dataset. As summarized in Table 1, the number of open-source whisper datasets is limited, and they exhibit significant variations in scale, linguistic diversity, and accessibility. The licensing terms (CC Type) for these datasets are also heterogeneous, including Creative Commons Attribution-NonCommercial-ShareAlike (NC-SA), Attribution-NoDerivatives (ND), Attribution-NonCommercial (NC), among others, while the license information for some datasets remains unspecified (N.A.). This heterogeneity directly impacts data availability, redistribution rights, and the freedom to utilize the data for academic and commercial research. In constructing **WhispReal**, we meticulously respected and adhered to the licensing agreements of all source datasets, ensuring the full legal compliance of the resulting collection.

Existing whispered speech corpora documented in the literature can be broadly categorized into three primary types, with a subset being publicly available as highlighted in gray in Table 1:

Dataset	Source	Sample Rate	Size (h)	Avg (s)	Language	Speaker	P?	CC Type
UTVE-I	(Zhang and Hansen, 2010)	44100	<1	N.A.	en	12	N	N.A.
UTVE-II	(Ghaffarzadegan et al., 2014)	44100	1	N.A.	en	112	N	N.A.
AVWD	(Zhou et al., 2019)	44100	<2.44	N.A.	zh	10	Y	N.A.
Whi-spe	(Grozdić et al., 2012)	22050	<5	N.A.	sr	10	Y	N.A.
AV-Whisper	(Tran et al., 2013)	48000	<10	N.A.	en	11	Y	N.A.
CIAIR	(Kawaguchi et al., 2002)	16000	15	N.A.	ja	123	Y	N.A.
iWhisper-Mandarin	(Lee et al., 2014a)	16000	15	N.A.	zh	80	Y	N.A.
wSPIRE	(Singhal et al., 2021)	44100	18	N.A.	en	88	Y	N.A.
AISHELL6-Whisper	(Li et al., 2025)	48000	29.75	9.01	zh	167	Y	BY-NC-SA
CHAINS (subset)	(Cummins et al., 2006)	44100	2.55	4.12	en	36	Y	BY-ND
EARs (subset)	(Richter et al., 2024)	48000	3.22	16.58	en	107	Y	BY-NC
Expresso (subset)	(Nguyen et al., 2023)	48000	1.32	3.13	en	4	Y	BY-NC
Whisper40	(Yang and Zhou, 2024)	16000	6.10	12.25	zh	40	Y	N.A.
wTIMIT	(Li-Li et al., 2005)	44100	29.44	5.04	en	48	Y	N.A.
WhispXXX	Our new curation	44100	45.24	9.01	zh	77	Y	BY-NC
WhispReal	Our aggregated corpus	Mixed	117.62	9.91	en, zh	479	Y	Mixed
WhispSynth	Synthesized from WhispReal	24000	≈118	≈10	en, zh	479	N	MIT

Table 1: Overview of open-source and newly introduced datasets used in this work, with the highlighted (gray) rows indicating all obtainable real whisper corpora included in our study. Entries include data duration (Size), average utterance length (Avg), language, speaker count, whether the paired normal segments are available (P?), and license type (CC Type), “Mixed” indicates multiple license terms or sample rate, “N.A.” denotes unknown or unspecified.

Pure Whisper Datasets These collections are specifically dedicated to whispered vocalizations, designed for targeted acoustic analysis (Li-Li et al., 2005; Jou et al., 2005; Zhang and Hansen, 2010; Lim, 2011; Ghaffarzadegan et al., 2014; Lee et al., 2014b; Singhal et al., 2021; Hiraki and Rekimoto, 2022; Yang and Zhou, 2024). Among these, wTIMIT (Li-Li et al., 2005) is the most notable and widely adopted. It is a classic, relatively large-scale English whispered speech corpus that also includes parallel normal speech. The linguistic content of these collections is largely confined to the word or phrase level.

Subsets within Expressive Speech Databases Several speech corpora designed for emotional or expressive analysis incorporate whispered speech as one vocalization style (Cummins et al., 2006; Nguyen et al., 2023; Richter et al., 2024). These datasets often contain rich emotional annotations, facilitating research on affective whispering, although the whisper portions themselves are typically limited in size.

Audiovisual Whisper Collections This is an emerging category that integrates whispered speech with synchronized visual recordings (Kawaguchi et al., 2002; Tran et al., 2013; Petridis et al., 2018; Zhou et al., 2019; Li et al., 2025). Within this domain, AISHELL6-Whisper (Li et al., 2025) represents one of the largest Chinese datasets, characterized by a substantial number of speakers, which makes it well-suited for multimodal analysis and robustness research. It is worth noting that the upcoming dataset, WhispXXX (anonymized), also falls

into this category. However, due to the heightened complexity of the recording process, the scalability of such audiovisual collections is constrained.

In the following sections, we present information on the existing available whispered data sources and a unique processing required for each. This is followed by a comprehensive account of the data collection procedures for our dataset.

2.1 Existing Available Resources

AISHELL6-Whisper (Li et al., 2025): it provides audio-visual whisper speech dataset, featuring 30 hours each of whisper speech and parallel normal speech, with synchronized frontal facial videos. We used only the audio and followed the train/test/valid splits. It should be noted that some utterances within the retained set still lack paired counterparts. **wTIMIT** (Lim, 2011): it is constructed after the TIMIT dataset (Garofolo et al., 1993). The first phase of the project collected the speech from 20 Singaporeans, and the second phase collected the speech from 28 North Americans. After excluding files with evident sampling errors, we retained all remaining audio files. It should be noted that some utterances within the retained set still lack whispered counterparts. Subsequently, we reorganized the dataset by speaker into a 60%–20%–20% split for training, validation, and testing, respectively.

Whisper40 (Yang and Zhou, 2024): it is a Mandarin Chinese corpus comprising whispered and normal speech recordings from 40 speakers without corresponding transcripts. However, since the normal and whispered speech pairs share identi-

cal linguistic content, the normal speech is clearly articulated, and the text is consistent across speakers, the authors performed manual transcription and verification for all utterances. We followed the train/test/valid splits of the original dataset.

CHAINS (Cummins et al., 2006): **CH**Aracterizing **IN**dividual **S**peakers contains recordings from 36 speakers (8 from the UK/US and 28 from Ireland), each producing 37 utterances across six different speaking conditions: solo speech, retelling, synchronous speech, repetitive synchronous imitation, fast speech, and whispered speech. We extracted the solo speech and whispered speech from CHAINS and segmented its four short fables (Cinderella, Rainbow text, North Wind and the Sun, and Members of the Body) into sentences to ensure that each audio file is under 30 seconds in duration, and split the data by speaker into 60% training, 20% validation, and 20% test sets.

Espresso (Nguyen et al., 2023): it is a high-quality (48 kHz) expressive speech dataset containing expressive read speech in 8 styles and improvised dialogues in 26 styles, recorded from 4 speakers. We selected default style and whisper style of it. Following a consistent partitioning scheme, we split the data by speaker into 60% training, 20% validation, and 20% test sets.

EARs (Richter et al., 2024) : **E**xpressive **A**nechoic **R**ecordings of **S**peech is a high-quality expressive speech collection, recorded at 48 kHz and comprising 107 speakers from diverse backgrounds, totaling approximately 100 hours of clean, anechoic speech. EARs covers the full dynamic range of vocal expression, from whispering to yelling and screaming. For this study, we extracted the whispered and regular speech subsets (split: 60% train / 20% valid / 20% test by speaker).

2.2 The WhispXXX corpus

The WhispXXX corpus consists of 85 hours of paired whispered and normal speech from 77 speakers of students recruited at X University. WhispXXX is designed to facilitate research and development in Mandarin whispered ASR. It is constructed after the THCHS-30 dataset (Wang and Zhang, 2015), which is usually used for the study of Mandarin speech recognition.

Speaker Recruitment The dataset was recorded by 77 native Mandarin speakers (37 male and 40 female), aged 22-26, recruited from X University. All participants met the following selection criteria essential for speech production research:

- Native proficiency in Mandarin Chinese.
- No diagnosed speech or hearing disorders.
- No acute oral/auditory health conditions.

These criteria ensured consistent recording quality while complying with ethical research standards. Participants received compensation and provided informed consent.

Record Settings The participants were instructed to produce each sentence twice: once in normal speech and once in whispered speech. Due to the time-consuming nature of the recordings, sessions could be divided into multiple sittings, and submissions were only accepted after all recordings were completed. To ensure universal applicability, no restrictions were placed on the speakers' regional origin or accent. During recording, no obstructions were allowed between the speaker and the recording device to avoid interfering with airflow pickup by the microphone.

Data preparation Following the partitioning scheme of the THCHS-30 dataset (Wang and Zhang, 2015), our dataset is divided into four groups according to the recording text: A (sentence ID from 1 to 250), B (sentence ID from 251 to 500), C (sentence ID from 501 to 750), and D(sentence ID from 751 to 1000). Each participant recorded 250 sentences (approximately 50 words per sentence), totaling approximately one hour of recording per speaker. Speakers were randomly assigned unique identification numbers and evenly distributed across the four groups.

3 WhispSynth

3.1 Generation Pipeline

Although recent advances in TTS synthesis have enabled the generation of highly natural and intelligible speech, when synthesizing whispered speech many state-of-the-art systems tend to produce residual pitch or fundamental frequency (F0) components. To address this issue, we propose a post-processing pipeline that detects and removes pitch components from synthesized whispers while preserving the desired noise-like whisper quality. Figure 4 outlines the proposed procedure.

The core idea is to leverage a DDSP decomposition to isolate and suppress harmonic content in pitch-contaminated segments, followed by an overlap-add (OLA) reconstruction to ensure smooth transitions. The full algorithm is provided in Appendix A.

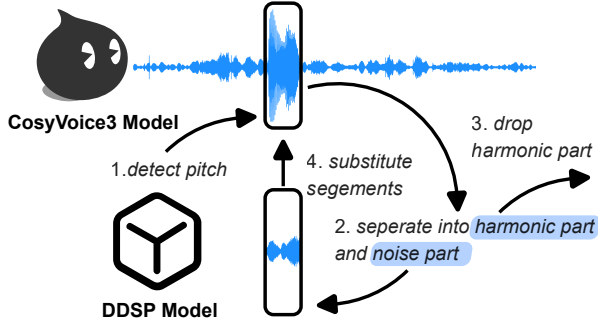


Figure 2: Visualization of the WhisperSynth’s generation pipeline. We apply CosyVoice3 and DDSP Model to generate pitch-free voice.

3.2 CosyVoice3: State-of-the-Art Speech Synthesis

CosyVoice3 (Du et al., 2025), developed by Alibaba Group, represents a next-generation speech synthesis model designed for high-quality and natural speech in open-environment scenarios. It integrates a multi-task speech tokenizer, scalable Differentiable Reward Optimization for post-training, and is trained on an ultra-large-scale multilingual corpus with instruction-following controllability.

The model exhibits particular strength in whisper synthesis. Its multi-task tokenizer, trained jointly using speech emotion recognition and audio event detection tasks, effectively encodes paralinguistic features such as softness and breathiness. Moreover, its million-hour-scale training corpus, drawn from diverse real-world recordings, contains abundant whisper and low-volume speech samples, enabling robust learning of such vocal qualities. Nonetheless, during whisper synthesis, the model occasionally exhibits fundamental frequency artifacts. This issue stems from the fact that the neural modules are predominantly trained on phonated speech, leading them to inadvertently introduce subtle periodicity even when generating whisper-like outputs.

3.3 Pitch-free with DDSP: An Overview

The DDSP model (Engel et al., 2020), proposed by researchers from Google, adopts the harmonic-plus-noise model (Serra and Smith, 1990) to decompose a monophonic sound into a harmonic component \mathbf{H} and a noise component \mathbf{N} :

$$\mathbf{S}[i] = \mathbf{H}[i] + \mathbf{N}[i]. \quad (1)$$

Given the upsampled estimated $\tilde{F}_0[i]$, in this work we approximate \mathbf{H} by a sawtooth signal (Wu et al.,

2022), which contains an equal number of even and odd harmonics with decaying magnitudes:

$$\tilde{\mathbf{H}}[i] = \sum_{k=1}^K \frac{1}{k} \sin(\phi_k[i]). \quad (2)$$

The instantaneous phase $\phi_k(t)$ is obtained by cumulative summation of the instantaneous frequency samples $k\tilde{F}_0[i]$:

$$\phi_k[i] = 2\pi \sum_{n=0}^N k\tilde{F}_0[i], \quad (3)$$

where $\tilde{\mathbf{H}}$ is treated as the “excitation signal” and shaped into the desirable \mathbf{y}_h by means of an linear time-varying finite impulse response (LTV-FIR) filter $\psi_h(i) \in R^{L_h}$:

$$\bar{\mathbf{H}}[i] = \tilde{\mathbf{H}}[i] * \psi_h[i], \quad (4)$$

where \mathbf{N} is approximated by convolving a uniform distributed noise signal ζ ranging from -1 to 1 (with the same length as a frame) with an LTV-FIR filter $\psi_n[i] \in R^{L_n}$ estimated per frame:

$$\bar{\mathbf{N}}[i] = \zeta[i] * \psi_n[i]. \quad (5)$$

Jointly, the parameters $\Phi := \{\tilde{F}_0[i], \psi_h[i], \psi_n[i]\}_{i=1}^N$ are estimated from an input mel-spectrogram \mathbf{X} per frame by an NN mapping function f_{NN} (Gulati et al., 2020):

$$\Phi = f_{\text{NN}}(\mathbf{X}). \quad (6)$$

DDSP entails a source-filter model (Wang et al., 2019) that can successfully separate the noise component from normal speech.

3.4 Pitch-free with DDSP: Training Pipeline

Based on the prior work of DDSP in singing vocoders (Wu et al., 2022), the authors identified a critical limitation: the model produces buzzing artifacts in unvoiced and semi-voiced segments when used as a subtractive synthesizer². To overcome this limitation and adapt DDSP for high-quality whispered speech synthesis, we propose the following three training strategies:

(i) **Adversarial Training with Normal Speech:** We integrate the DDSP generator into the BigVGAN (Lee et al., 2022) framework, applying a multi-resolution discriminator (MRD) and a

²<https://github.com/YatingMusic/ddsp-singing-vocoders/tree/main/postprocessing>

multi-period discriminator (SPD) with least-square GAN loss \mathcal{L}_{adv} (Mao et al., 2017). This adversarial training is essential because GAN-based discrimination has been shown to effectively improve the modeling of harmonic structures and suppress artifact-prone synthesis patterns.

(ii) **Continued Training with Whispered Speech:** After initially training the DDSF vocoder on normal speech data, we continue training it on the WhispReal dataset without adversarial objectives for faster training speed. This two-stage procedure allows the model to first learn robust pitch-conditioned synthesis, then specialize in whispered speech generation while avoiding instability that might arise from direct training on limited whispered data.

(iii) **Semi-supervised Dual-focus Training:** We observe that the WhispReal dataset contains samples with perceptible pitch contours resembling normal speech, due to speaker errors or temporary vocal fatigue. Therefore the training scheme always treats the output as a sum of harmonic and stochastic components, but shifts the learning emphasis depending on the input type: whispered/normal speech share loss, but the gradient reweights harmonic/noise terms so each half teaches the other. This dual-focus approach enables mutually reinforcing improvements in both harmonic and stochastic modeling.

4 CosyWhisper

We fine-tune CosyVoice3 with Whisper, following the exact procedure outlined in the official script. The CosyVoice3 architecture comprises three primary components: the text-to-speech large language model (LLM), the Conditional Flow Matching (CFM) model, and HiFi-GAN vocoder. In our scenario, the distinction between whispered and normal speech lies primarily in acoustic modeling rather than semantic content. Therefore, we selectively fine-tune only the CFM model, which is responsible for converting semantic tokens into high-fidelity acoustic features, while keeping the LM and HiFi-GAN unchanged.

Although the official script currently supports only LLM training (line 65: “Run train. We only support LLM training for now”), we successfully extend it to CFM fine-tuning. Key adaptations include adding a token projection layer, replacing the encoder with a direct embedding lookup, and revising the conditioning mechanism. These changes

improve token-acoustic alignment and model efficiency. Implementation specifics are provided in the Appendix B. We refer to the resulting fine-tuned model as **CosyWhisper**, a name that reflects its origin in the CosyVoice3 architecture and its specialized capability in Whisper-based speech synthesis.

5 Experiments

5.1 Experimental setups

Data Settings We used the English and Mandarin Chinese subsets of the WhispReal dataset for training and evaluation. And the WhispSynth dataset is partitioned accordingly. For each source, data were divided as described in Section 2.1. The validation set was used for model checkpoint selection, while the final sound quality evaluation was performed only on the test set. Detailed statistics of the dataset distribution are provided in Table 2.

	en		zh	
	# Train	# Test	# Train	# Test
File Count	20694	4770	32704	7256
Size (h)	29.62	6.91	65.86	15.25
Avg (s)	5.15	5.21	7.25	7.57
Speaker Count	78M 87F	15M 15F	206M 194F	26M 25F

Table 2: The statistics of the WhispReal dataset train/test distribution.

Training Settings We transformed the normalized speech in all the datasets into mel-spectrograms with a frame length of 1280 and a hop length of 320. Adversarial training was conducted using a batch size of 8, while standard (non-adversarial) training employed a batch size of 32. All experiments were performed on eight V100-32GB GPUs.

Evaluation Settings The generation performance is evaluated using both subjective and objective metrics. In terms of subjective evaluation, the proposed and baseline systems were evaluated using a 5-point Whisper-likeness Mean Opinion Score (W-MOS). It measured the perceptual similarity between the synthesized whisper and the listener’s mental prototype of a whisper, ranging from 0 (“Totally Not similar”) to 5 (“Extremely similar”). We have developed an automated subjective evaluation interface, as shown in Appendix C. For objective assessment, we adopt standard TTS metrics, as no widely used metrics are specifically designed for whisper. For the naturalness, we measure spectral differences with Mel Cepstral Distortion (MCD) (Chen et al., 2022), DNSMOS (Reddy et al., 2021), and UTMOS (Saeki et al., 2022)

Dataset	Source	Language	Naturalness \uparrow		Intelligibility \downarrow	F0 \downarrow
			DNSMOS	UTMOS	CER / WER (%)	VTR(%)
wTIMIT	(Li-Li et al., 2005)	en	2.76	1.31	50.99	0.69
CHAINS (<i>subset</i>)	(Cummins et al., 2006)	en	2.76	1.48	11.55	0.95
Espresso (<i>subset</i>)	(Nguyen et al., 2023)	en	3.24	1.50	7.77	0.58
EARs (<i>subset</i>)	(Richter et al., 2024)	en	3.43	2.01	6.31	0.30
Whisper40	(Yang and Zhou, 2024)	zh	2.61	1.27	72.80	0.99
AISHELL6-Whisper	(Li et al., 2025)	zh	2.75	1.68	38.26	0.94
WhispXXX	Ours	zh	2.90	1.33	36.74	0.98
WhispReal	Ours	en,zh	2.80	1.44	39.30/37.58	0.88
WhispSynth	Ours	en,zh	2.89	1.46	31.16/20.98	0.87

Table 3: Objective evaluation results on open-sourced and our proposed datasets. Comprehensive evaluation on whisper speech conversion and synthesis. **Quality**: Subjective naturalness (DNSMOS/UTMOS \uparrow). **Intelligibility**: Character/Word Error Rate (\downarrow). **Pitch**: Voiced Time Ratio (VTR) (\downarrow).

to estimate perceptual audio quality. We likewise calculate the Word Error Rate (WER) (Wang et al., 2018) for English and Character Error Rate (CER) (Xu et al., 2025) for Chinese to gauge intelligibility using SOTA multilingual ASR model Fun-ASR-nano-2512 (An et al., 2025). For the timbre, we calculate cosine similarity metrics based on ECAPA-TDNN (Desplanques et al., 2020) to obtain speaker identity similarity (SpkSim). Additionally, we evaluate Voiced Time Ratio (VTR) for pitch-free assessment.

5.2 Evaluating on Benchmarks

Based on the objective evaluation results in Table 3, several key observations can be made regarding the quality of whispered corpora. Among the publicly available English datasets, the EARs subset demonstrates superior performance, achieving the highest naturalness scores (DNSMOS: 3.43, UTMOS: 2.01) and the lowest CER: 6.31%, indicating its high-quality whisper characteristics. For Chinese datasets, our proposed WhispXXX corpus shows a favorable balance, attaining a higher DNSMOS score (2.90) compared to other Chinese whisper collections like Whisper40 (2.61) and AISHELL6-Whisper (2.75). The primary advantages of our newly introduced datasets, WhispReal and WhispSynth, are threefold. First, they provide multilingual coverage for both English and Chinese, addressing a gap in existing resources that are predominantly monolingual. Second, WhispSynth exhibits exceptional intelligibility, achieving the lowest error rates (CER: 31.16%, WER: 20.98%) among all evaluated datasets, which is beneficial for speech recognition and synthesis tasks requiring high clarity. Third, it maintains a balanced naturalness, with competitive DNSMOS (2.89) and UTMOS (1.46) scores while demonstrating a reasonable pitch profile (VTR: 0.87) comparable to

high-quality datasets like Espresso.

Overall, WhispSynth presents an optimal trade-off between intelligibility and naturalness, making it a suitable resource for whisper-based speech processing applications. In addition, the results also reveal that existing objective metrics are poorly aligned with whisper characteristics; ultimately, human judgment remains indispensable.

5.3 Experiments for generated samples

5.3.1 Comparative Systems

Whisper-Effect (Roh et al., 2025) is a designed acoustic perturbation that simulates whisper by applying a high-pass filter combined with white noise. **toWhisper** (zeta, 2017) is an Linear Predictive Coding (LPC)-based tool that synthesizes whispered speech by processing white noise through a vocal tract filter estimated via LPC vocoding.

Normal2Whisper (Lin et al., 2023) converts normal speech to whisper through a two-stage process (Glottal source removal and Spectral modification) by the WORLD vocoder (Morise et al., 2016).

SeedVC (Liu, 2024) is a zero-shot voice conversion framework that disentangles representations of content and speaker information.

CosyVoice3 (Du et al., 2025) is a multilingual TTS model for high-quality and efficient generation.

5.4 Experiments for CosyWhisper

5.4.1 Quantitative Evaluation

The objective results in Table 4 highlight two key findings. First, in *Normal-to-Whisper Conversion*, voice conversion models like SeedVC yield high speaker similarity and intelligibility, but their whisperiness remains limited (W-MOS: 2.18). Although traditional baselines (e.g., toWhisper) are non-trainable, lack adaptability, and yield lower W-MOS despite decent objective scores. Our pitch-

Method	Source	Whisperiness		Naturalness		Intelligibility	Timbre	Pitch
		W-MOS \uparrow	DNSMOS \uparrow	UTMOS \uparrow	MCD \downarrow	CER/WER (%) \downarrow	SpkSim \uparrow	VTR (%) \downarrow
Ground Truth	Test Set of WhispReal	4.33 \pm 0.33	2.80	1.44	0.00	39.30/37.58	1.00	0.88
<i>Normal-to-Whisper Conversion</i>								
Whisper-Effect toWhisper	(Roh et al., 2025) (zeta, 2017)	-	3.14	1.30	70.23	59.46/80.35	0.52	0.12
Normal2Whisper	(Lin et al., 2023)	1.02 \pm 0.29	2.94	1.35	55.73	12.07/ 9.27	0.58	0.92
SeedVC	(Liu, 2024)	-	2.89	1.28	60.00	11.92/28.28	0.59	0.94
Pitch-free Model	Ours	2.18 \pm 0.47	3.01	1.61	54.85	58.65/19.97	0.66	0.71
Pitch-free Model	Ours	1.30 \pm 0.29	2.90	1.36	63.05	21.31/45.18	0.63	0.98
<i>Text-to-Whisper Synthesis</i>								
CosyVoice3	(Du et al., 2025)	3.40 \pm 0.51	3.00	1.47	56.03	12.51/9.81	0.83	0.86
CosyWhisper	Ours	4.53\pm0.20	3.08	1.48	59.31	12.76/29.22	0.80	0.88

Table 4: Evaluation on whisper speech conversion and synthesis. **Quality**: Objective whisperiness (W-MOS \uparrow), Subjective naturalness (DNSMOS/UTMOS \uparrow) and objective distortion (MCD \downarrow). **Intelligibility**: Character/Word Error Rate (\downarrow). **Timbre**: Speaker similarity cosine similarity (\uparrow). **Pitch**: Voiced Time Ratio (VTR) (\downarrow).

free model is trainable and better captures natural whisper variability, confirming that not fixed transformations for authentic whisper synthesis. Second, in *Text-to-Whisper Synthesis*, even strong TTS systems like CosyVoice3 produce speech that is intelligible and natural, yet still perceptibly non-whisper (W-MOS: 3.40). In contrast, our CosyWhisper achieves a W-MOS of 4.53—the highest among all synthesized approaches and even better than the ground-truth (4.33), demonstrating our superiority to generate authentic whisper from text.

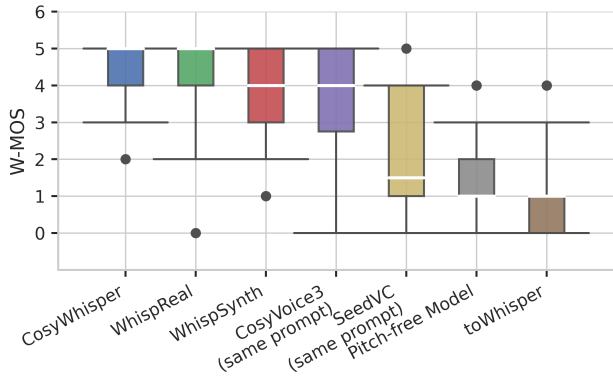


Figure 3: Results of the listening test. Whisper-likeness Mean Opinion Score (W-MOS) based on 20 participants visualized in a standard box plot.

For subjective evaluation, the W-MOS box plot from the listening test statistically validates that listeners perceive the output of CosyWhisper as similar to natural whisper. The interquartile range (IQR) for CosyWhisper is expected to be compact and positioned high on the scale (near 4.0), with a median close to 4.1, indicating consistent high-quality ratings from participants with low disagreement. In contrast, the boxes for other methods (e.g., CosyVoice3, SeedVC) would be located lower, with potentially larger IQRs, reflecting greater vari-

ability in listener perception and lower consensus on their whisper authenticity. CosyWhisper successfully bridges the gap between high intelligibility and authentic whisper perceptual quality.

5.5 Ablation on Tuning Dataset

To evaluate the efficacy of our synthetic data for text-to-whisper synthesis, we fine-tune our CosyWhisper on two variants of the corpus: WhispReal and our synthetic WhispSynth. As shown in Table 5, WhispSynth consistently outperforms across all key metrics: it reduces CER/WER by 46%, lowers pitch distortion (VTR) by 9%, and improves naturalness by 8%, despite being fully synthetic. This confirms that WhispSynth not only mitigates the noise and inconsistency inherent in real whisper data but also provides high-fidelity data that better supports stable whisper generation from text.

Dataset	Naturalness \uparrow		Intelligibility \downarrow	Timbre \uparrow	F0 \downarrow
	DNSMOS	UTMOS	CER / WER (%)	Cosine	VTR (%)
WhispReal	2.93	1.33	28.3/46.5	0.75	0.77
WhispSynth	3.08	1.48	12.8/29.2	0.80	0.70

Table 5: Performance comparison of our CosyWhisper fine-tuned on WhispReal versus WhispSynth.

6 Conclusion

In this paper, we address the fundamental data scarcity challenge in whisper research by introducing WhispSynth, a large-scale, high-fidelity multilingual whisper corpus. Our core contribution is a novel generative framework that integrates a DDSP-based pitch-free method with advanced TTS models, transforming diverse and noisy real whispered recordings into a clean, studio-grade synthetic dataset. This pipeline ensures the faithful preservation of vocal timbre and linguistic content while significantly enhancing acoustic quality.

617 **Limitations**

618 The impact of non-linguistic variations (e.g., the
619 use of different microphones) on model perfor-
620 mance was not systematically assessed. Although
621 existing research suggests that CosyVoice3 are rel-
622 atively robust to such variations, the influence of
623 hardware differences in real-world deployment sce-
624 narios on specific task performance requires further
625 verification. Besides, to address security concerns,
626 the CosyWhisper model used in this study will be
627 released with an automatically embedded real-time
628 audio watermark. While this measure is crucial for
629 responsible usage tracking, it may have a potential
630 impact on the acoustic properties of the synthesized
631 audio and subsequent analyses.

632 **Ethics Statement**

633 Our paper evaluated various methods that could
634 make developing text-to-whisper synthesize sys-
635 tems more viable for languages where paired whis-
636 per and transcriptions are difficult to obtain. In
637 our experiments, we only used already publicly
638 available data (CHAINS, EARS, Espresso, Whis-
639 per40) or data for which we have obtained informed
640 consent for public release from the data custodi-
641 ans (AISHELL6-Whisper, wTIMIT). To make our
642 findings as relevant as possible for other language
643 projects, we minimized the amount of computing
644 time used.

645 **References**

646 Keyu An, Yanni Chen, Chong Deng, Changfeng Gao,
647 Zhifu Gao, Bo Gong, Xiangang Li, Yabin Li, Xiang
648 Lv, Yunjie Ji, and 1 others. 2025. Fun-asr technical
649 report. *arXiv preprint arXiv:2509.12508*.

650 Joceline Andersen. 2015. Now you’ve got the shiveries:
651 Affect, intimacy, and the asmr whisper community.
652 *Television & New Media*, 16(8):683–700.

653 ByteDance. 2025. *Doubao ai chatbot*. <https://www.doubao.com/chat>. Accessed: [2025-07-01].

655 Qi Chen, Mingkui Tan, Yuankai Qi, Jiaqiu Zhou, Yuan-
656 qing Li, and Qi Wu. 2022. V2C: Visual voice cloning.
657 In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages
658 21210–21219.

659 Yushen Chen, Zhikang Niu, Ziyang Ma, Keqi Deng,
660 Chunhui Wang, JianZhao JianZhao, Kai Yu, and Xie
661 Chen. 2025. F5-tts: A fairytaler that fakes fluent and
662 faithful speech with flow matching. In *Proceedings*
663 *of the 63rd Annual Meeting of the Association for*
664 *Computational Linguistics (Volume 1: Long Papers)*,
665 pages 6255–6271.

Fred Cummins, Marco Grimaldi, Thomas Leonard, and
666 Juraj Simko. 2006. The chains corpus: Characteriz-
667 ing individual speakers. In *Proc of SPECOM*, pages
668 1–6. 669

Brecht Desplanques, Jenthe Thienpondt, and Kris De-
670 muynck. 2020. ECAPA-TDNN: emphasized chan-
671 nel attention, propagation and aggregation in TDNN
672 based speaker verification. In *Annu. Conf. Int. Speech*
673 *Commun. Assoc.*, pages 3830–3834. 674

Zhihao Du, Changfeng Gao, Yuxuan Wang, Fan
675 Yu, Tianyu Zhao, Hao Wang, Xiang Lv, Hui
676 Wang, Chongjia Ni, Xian Shi, and 1 others. 2025.
677 Cosyvoice 3: Towards in-the-wild speech genera-
678 tion via scaling-up and post-training. *arXiv preprint*
679 *arXiv:2505.17589*. 680

ElevenLabs. 2025. *Elevenlabs: Free text to speech &*
681 *ai voice generator*. <https://elevenlabs.io>. Ac-
682 cessed: [2025-07-01]. 683

Jesse Engel, Chenjie Gu, Adam Roberts, and 1 others.
684 2020. Ddsp: Differentiable digital signal processing.
685 In *International Conference on Learning Representa-*
686 *tions (ICLR)*. 687

John S Garofolo, Lori F Lamel, William M Fisher,
688 Jonathan G Fiscus, and David S Pallett. 1993. Darpa
689 timit acoustic-phonetic continuous speech corpus cd-
690 rom. nist speech disc 1-1.1. *NASA STI/Recon Techni-*
691 *cal Report N*, 93:27403. 692

Shabnam Ghaffarzagdegan, Hynek Bořil, and John HL
693 Hansen. 2014. Ut-vocal effort ii: Analysis and
694 constrained-lexicon recognition of whispered speech.
695 In *2014 IEEE International Conference on Acous-*
696 *tics, Speech and Signal Processing (ICASSP)*, pages
697 2544–2548. IEEE. 698

Đorđe T Grozdić, Branko Marković, Jovan Galić, and
699 Slobodan T Jovičić. 2012. Application of neural
700 networks in whispered speech recognition. In *2012*
701 *20th Telecommunications Forum (TELFOR)*, pages
702 728–731. IEEE. 703

Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki
704 Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang,
705 Zhengdong Zhang, Yonghui Wu, and 1 others. 2020.
706 Conformer: Convolution-augmented transformer for
707 speech recognition. In *Proc. Interspeech 2020*, pages
708 5036–5040. 709

Hirotaaka Hiraki and Jun Rekimoto. 2022. Silentwhis-
710 per: Faint whisper speech detection using wearable
711 microphones. In *Adjunct Proceedings of the 35th*
712 *Annual ACM Symposium on User Interface Software*
713 *and Technology*, pages 1–3. 714

Hirotaaka Hiraki and Jun Rekimoto. 2025. Silentwhis-
715 per: inaudible faint whisper speech input for silent
716 speech interaction. In *Proceedings of the Extended*
717 *Abstracts of the CHI Conference on Human Factors*
718 *in Computing Systems*, pages 1–6. 719

720	Taisuke Ito, Kazuya Takeda, and Fumitada Itakura.	Masanori Morise, Fumiya Yokomori, and Kenji Ozawa.	773
721	2005. Analysis and recognition of whispered speech.	2016. World: A vocoder-based high-quality speech	774
722	<i>Speech Communication</i> .	synthesis system for real-time applications. <i>IE-</i>	775
723	Szu-Chen Jou, Tanja Schultz, and Alex Waibel. 2005.	<i>ICE Transactions on Information and Systems</i> ,	776
724	Whispery speech recognition using adapted articu-	99(7):1877–1884.	777
725	latory features. In <i>2005 IEEE International Confer-</i>	Tu Anh Nguyen, Wei-Ning Hsu, Antony D’Avirro,	778
726	<i>ence on Acoustics, Speech, and Signal Processing</i>	Bowen Shi, Itai Gat, Maryam Fazel-Zarandi, Tal Re-	779
727	<i>(ICASSP)</i> , volume 1, pages 1009–1012. IEEE.	mez, Jade Copet, Gabriel Synnaeve, Michael Hassid,	780
728	N Kawaguchi, K Takeda, S Matsubara, I Yokoo, T Ito,	and 1 others. 2023. Expresso: A benchmark and	781
729	K Tatara, T Shinde, and F Itakura. 2002. Ciair speech	analysis of discrete expressive speech resynthesis. In	782
730	corpus for real world applications. In <i>The Interna-</i>	<i>Proc. Interspeech 2023</i> , pages 4823–4827.	783
731	<i>tional conference Committee for the Coordination</i>	Stavros Petridis, Jie Shen, Doruk Cetin, and Maja Pantic.	784
732	<i>and Standardization of Speech Databases and Asses-</i>	2018. Visual-only recognition of normal, whispered	785
733	<i>ment Techniques</i> .	and silent speech. In <i>2018 IEEE International Con-</i>	786
734	Pei Xuan Lee, Darren Wee, Hilary Si Yin Toh,	<i>ference on Acoustics, Speech and Signal Processing</i>	787
735	Boon Pang Lim, Nancy F Chen, and Bin Ma. 2014a.	<i>(ICASSP)</i> , pages 6219–6223. IEEE.	788
736	A whispered mandarin corpus for speech technology	Chandan K. A. Reddy, Vishak Gopal, and Ross Cutler.	789
737	applications. In <i>INTERSPEECH</i> , pages 1598–1602.	2021. Dnsmos: A non-intrusive perceptual objective	790
738	Pei Xuan Lee, Darren Wee, Hilary Si Yin Toh,	speech quality metric to evaluate noise suppressors.	791
739	Boon Pang Lim, Nancy F Chen, and Bin Ma. 2014b.	In <i>IEEE Conf. Acoust. Speech Signal Process.</i> , pages	792
740	A whispered mandarin corpus for speech technol-	6493–6497.	793
741	ogy applications. In <i>Proc. Interspeech 2014</i> , pages	Jun Rekimoto. 2023. Wesper: Zero-shot and real-	794
742	1598–1602.	time whisper to normal voice conversion for whisper-	795
743	Sang-gil Lee, Wei Ping, Boris Ginsburg, Bryan Catan-	based speech interactions. In <i>Proceedings of the</i>	796
744	zaro, and Sungroh Yoon. 2022. Bigvgan: A univer-	<i>2023 CHI conference on human factors in computing</i>	797
745	sal neural vocoder with large-scale training. <i>Inter-</i>	<i>systems</i> , pages 1–12.	798
746	<i>national Conference on Learning Representations</i>	Julius Richter, Yi-Chiao Wu, Steven Krenn, Simon	799
747	<i>(ICLR)</i> .	Welker, Bunlong Lay, Shinji Watanabe, Alexander	800
748	Cancan Li, Fei Su, Juan Liu, Hui Bu, Yulong Wan,	Richard, and Timo Gerkmann. 2024. Ears: An	801
749	Hongbin Suo, and Ming Li. 2025. Aishell6-whisper:	anechoic fullband speech dataset benchmarked for	802
750	A chinese mandarin audio-visual whisper speech	speech enhancement and dereverberation. In <i>Proc.</i>	803
751	dataset with speech recognition baselines. <i>arXiv</i>	<i>Interspeech 2024</i> , pages 4873–4877.	804
752	<i>preprint arXiv:2509.23833</i> .	Jaechul Roh, Virat Shejwalkar, and Amir Houmansadr.	805
753	Yang Li-Li, Li Yan, and Xu Bo-Ling. 2005. The es-	2025. Multilingual and multi-accent jailbreaking of	806
754	tablishment of a chinese whisper database and per-	audio llms. <i>arXiv preprint arXiv:2504.01094</i> .	807
755	ceptual experiment. <i>Journal of Nanjing University</i>	Takaaki Saeki, Detai Xin, Wataru Nakata, Tomoki	808
756	<i>(Natural Sciences Edition)</i> , 41(3):311.	Koriyama, Shinnosuke Takamichi, and Hiroshi	809
757	Boon Pang Lim. 2011. <i>Computational Differences Be-</i>	Saruwatari. 2022. Utmos: Utokyo-sarulab sys-	810
758	<i>tween Whispered and Non-whispered Speech</i> . Uni-	tem for voicemos challenge 2022. <i>arXiv preprint</i>	811
759	versity of Illinois at Urbana-Champaign.	<i>arXiv:2204.02152</i> .	812
760	Zhaofeng Lin, Tanvina Patel, and Odette Scharenborg.	Xavier Serra and Julius Smith. 1990. Spectral mod-	813
761	2023. Improving whispered speech recognition per-	eling synthesis: A sound analysis/synthesis system	814
762	formance using pseudo-whisper based data augmen-	based on a deterministic plus stochastic decomposi-	815
763	tation. In <i>IEEE Automatic Speech Recognition and</i>	tion. <i>Computer Music Journal</i> , 14(4):12–24.	816
764	<i>Understanding Workshop (ASRU)</i> , pages 1–8. IEEE.	Bhavuk Singhal, Abinay Reddy Naini, and Prasanta Ku-	817
765	Songting Liu. 2024. Zero-shot voice conversion	mar Ghosh. 2021. wspire: A parallel multi-device	818
766	with diffusion transformers. <i>arXiv preprint</i>	corpus in neutral and whispered speech. In <i>2021 24th</i>	819
767	<i>arXiv:2411.09943</i> .	<i>Conference of the Oriental COCOSDA International</i>	820
768	Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau,	<i>Committee for the Co-ordination and Standardisa-</i>	821
769	Zhen Wang, and Stephen Paul Smolley. 2017. Least	<i>tion of Speech Databases and Assessment Techniques</i>	822
770	squares generative adversarial networks. In <i>Proceed-</i>	<i>(O-COCOSDA)</i> , pages 146–151. IEEE.	823
771	<i>ings of the IEEE International Conference on Com-</i>	VC Tartter. 1989. What’s in a whisper? <i>The Journal of</i>	824
772	<i>puter Vision (ICCV)</i> , pages 2794–2802.	<i>the Acoustical Society of America</i> , 86(5):1678–1683.	825

826 Tam Tran, Soroosh Mariooryad, and Carlos Busso. 2013.
827 Audiovisual corpus to analyze whisper speech. In
828 *2013 IEEE International Conference on Acoustics,
829 Speech and Signal Processing (ICASSP)*, pages 8101–
830 8105. IEEE.

831 Dong Wang and Xuewei Zhang. 2015. Thchs-30:
832 A free chinese speech corpus. *arXiv preprint
833 arXiv:1512.01882*.

834 Xin Wang, Shinji Takaki, and Junichi Yamagishi. 2019.
835 Neural source-filter waveform models for statistical
836 parametric speech synthesis. *IEEE/ACM Transac-
837 tions on Audio, Speech, and Language Processing*,
838 28:402–415.

839 Yuxuan Wang, Daisy Stanton, Yu Zhang, R. J. Skerry-
840 Ryan, Eric Battenberg, Joel Shor, Ying Xiao, Ye Jia,
841 Fei Ren, and Rif A. Saurous. 2018. Style tokens:
842 Unsupervised style modeling, control and transfer in
843 end-to-end speech synthesis. In *Int. Conf. on Mach.
844 Learn.*, volume 80, pages 5167–5176.

845 Da-Yi Wu, Wen-Yi Hsiao, Fu-Rong Yang, Oscar Fried-
846 man, Warren Jackson, Scott Bruzenak, Yi-Wen Liu,
847 and Yi-Hsuan Yang. 2022. Ddsp-based singing
848 vocoders: A new subtractive-based synthesizer and
849 comprehensive evaluation. In *International Society
850 for Music Information Retrieval Conference (ISMIR)*.

851 Kai-Tuo Xu, Feng-Long Xie, Xu Tang, and Yao
852 Hu. 2025. Fireredasr: Open-source industrial-
853 grade mandarin speech recognition models from
854 encoder-decoder to llm integration. *arXiv preprint
855 arXiv:2501.14350*.

856 Jingwen Yang and Ruohua Zhou. 2024. Whisper40:
857 A multi-person chinese whisper speaker recognition
858 dataset containing same-text neutral speech. *Infor-
859 mation*, 15(4):184.

860 zeta. 2017. *towhisper*. [https://github.com/
861 zeta-chicken/toWhisper](https://github.com/zeta-chicken/toWhisper). Accessed: [2025-07-
862 01].

863 Chi Zhang and John HL Hansen. 2010. Whisper-
864 island detection based on unsupervised segmentation
865 with entropy-based speech feature processing. *IEEE
866 Transactions on Audio, Speech, and Language Pro-
867 cessing*, 19(4):883–894.

868 Jian Zhou, Yuting Hu, Hailun Lian, Cong Pang, Huabin
869 Wang, and Liang Tao. 2019. An audio-visual whisper
870 database in chinese. In *Journal of Physics: Confer-
871 ence Series*, volume 1237, page 022106. IOP Pub-
872 lishing.

873 Siyi Zhou, Yiquan Zhou, Yi He, Xun Zhou, Jinchao
874 Wang, Wei Deng, and Jingchen Shu. 2025. In-
875 dexotts2: A breakthrough in emotionally expressive
876 and duration-controlled auto-regressive zero-shot
877 text-to-speech. *arXiv preprint arXiv:2506.21619*.

878 Appendix

879 A Algorithm

Algorithm 1 outlines the proposed procedure.

Algorithm 1: Whisper Synthesis with Pitch-Aware Segment Replacement

Input: Prompt audio \mathbf{A}_p , Corresponding text T_p , Instruction I
Output: Synthesized whisper audio \mathbf{A}_w . Initialize $\mathbf{A}_w \leftarrow \emptyset$ // Final whisper audio
 $\mathbf{A}_w^{\text{init}} \leftarrow \text{CosyVoice3}(\mathbf{A}_p, T_p, I)$
 $\mathcal{P} \leftarrow \text{DDSP_PitchDetection}(\mathbf{A}_w^{\text{init}})$
if $\mathcal{P} \neq \emptyset$ **then**
 foreach pitch segment $\mathbf{S}_i \in \mathcal{P}$ **do**
 $(\mathbf{H}_i, \mathbf{N}_i) \leftarrow \text{DDSP_Decompose}(\mathbf{S}_i)$
 // \mathbf{H}_i : harmonic segment, \mathbf{N}_i : noise segment
 $\mathbf{S}_i^{\text{noise}} \leftarrow \mathbf{N}_i$
 $\mathbf{S}_i^{\text{replaced}} \leftarrow \text{ApplyWindow}(\mathbf{S}_i^{\text{noise}})$
 $\mathbf{A}_w^{\text{init}} \leftarrow \text{ReplaceSegment}(\mathbf{A}_w^{\text{init}}, \mathbf{S}_i, \mathbf{S}_i^{\text{replaced}})$
 $\mathbf{A}_w \leftarrow \text{OverlapAdd}(\mathbf{A}_w^{\text{init}})$
else
 $\mathbf{A}_w \leftarrow \mathbf{A}_w^{\text{init}}$ // No pitch detected, use original
return \mathbf{A}_w

880 First, an initial whisper is generated using
881 CosyVoice3, conditioned on a prompt audio, its
882 transcript; note that a instruction is optional. The
883 pitch detection module in DDSP then identifies
884 segments containing residual F0. For each such
885 segment, DDSP decomposition separates the har-
886 monic and noise components. The harmonic part
887 is dropped, while the noise component, which car-
888 ries the whisper’s spectral envelope and aperiodic
889 excitation, is retained. Finally, an OLA opera-
890 tion is applied to reconstruct a continuous, pitch-free
891 whisper waveform.
892

893 B Code

894 Below are two versions of the code: the original
895 implementation (Original) and the revised version
896 (Modified). After implementing these modifica-
897 tions in the file³, we were able to successfully fine-
898 tune the Flow model within CosyVoice3.
899

```

899 def forward(
900     self,
901     batch: dict,
902     device: torch.device,
903 ) -> Dict[str, Optional[torch.Tensor]]:
904     # Data extraction
905
```

³<https://github.com/FunAudioLLM/CosyVoice/blob/652132ebaa3133428bf4db7e092a3cd1e073ca80/examples/libritts/cosyvoice3/cosyvoice/flow/flow.py>

```

906 8 token = batch['speech_token'].to(device)
907 9 token_len = batch['speech_token_len'].to(device)
908 10 feat = batch['speech_feat'].to(device)
909 11 feat_len = batch['speech_feat_len'].to(device)
910 12 embedding = batch['embedding'].to(device)
911 13
912 14 # NOTE unified training, static_chunk_size > 0 or = 0
913 15 streaming = True if random.random() < 0.5 else False
914 16
915 17 # xvec projection
916 18 embedding = F.normalize(embedding, dim=1)
917 19 embedding = self.spk_embed_affine_layer(embedding)
918 20
919 21 # concat text and prompt_text
920 22 mask = (
921 23     (~make_pad_mask(token_len))
922 24     .float()
923 25     .unsqueeze(-1)
924 26     .to(device)
925 27 )
926 28
927 29 token = (
928 30     self.input_embedding(torch.clamp(token, min=0))
929 31     * mask
930 32 )
931 33
932 34 # text encode
933 35 h, h_lengths = self.encoder(
934 36     token,
935 37     token_len,
936 38     streaming=streaming
937 39 )
938 40 h = self.encoder_proj(h)
939 41
940 42 # get conditions
941 43 conds = torch.zeros(feat.shape, device=token.device)
942 44 for i, j in enumerate(feat_len):
943 45     if random.random() < 0.5:
944 46         continue
945 47         index = random.randint(0, int(0.3 * j))
946 48         conds[i, :index] = feat[i, :index]
947 49 conds = conds.transpose(1, 2)
948 50
949 51 mask = (
950 52     ~make_pad_mask(
951 53         h_lengths.sum(dim=-1).squeeze(dim=1)
952 54     )
953 55     .to(h)
954 56 )
955 57 loss, _ = self.decoder.compute_loss(
956 58     feat.transpose(1, 2).contiguous(),
957 59     mask.unsqueeze(1),
958 60     h.transpose(1, 2).contiguous(),
959 61     embedding,
960 62     cond=conds,
961 63     streaming=streaming,
962 64 )
963 65 return {'loss': loss}

```

Listing 1: Original

```

964 1 self.token_proj = nn.Linear(input_size, input_size)
965 2
966 3 def forward(
967 4     self,
968 5     batch: dict,
969 6     device: torch.device,
970 7 ) -> Dict[str, Optional[torch.Tensor]]:
971 8
972 9     # 1. Data extraction
973 10 token = batch['speech_token'].to(device)
974 11 token_len = batch['speech_token_len'].to(device)
975 12 feat = batch['speech_feat'].to(device)
976 13 feat_len = batch['speech_feat_len'].to(device)
977 14 embedding = batch['embedding'].to(device)
978 15
979 16
980 17 streaming = random.random() < 0.5
981 18
982 19 # 2. Speaker embedding processing
983 20 embedding = F.normalize(embedding, dim=-1)
984 21 embedding = self.spk_embed_affine_layer(embedding)
985 22
986 23 # 3. Token to embedding mapping
987 24 # Ensure token shape is [B, T]
988 25 if token.dim() == 3:
989 26     token = token.squeeze(-1)
990 27
991 28 token = token.long()
992 29
993 30 # Token mask: [B, T, 1]
994 31 token_mask = (
995 32     (~make_pad_mask(token_len))
996 33     .unsqueeze(-1)
997 34     .to(device)
998 35 )
999 36
1000 37 # Embedding lookup

```

```

38 h = self.input_embedding(token) * token_mask
39 h_lengths = token_len
40
41 # 4. Token-to-mel alignment
42 h = h.repeat_interleave(self.token_mel_ratio, dim=1)
43 token_mask = token_mask.repeat_interleave(
44     self.token_mel_ratio,
45     dim=1
46 )
47 h_lengths = h_lengths * self.token_mel_ratio
48
49 # 5. Conditional mel (cond) preparation
50 conds = torch.zeros_like(feat)
51 for i, j in enumerate(feat_len):
52     if random.random() < 0.5:
53         continue
54         index = random.randint(0, int(0.3 * j))
55         conds[i, index] = feat[i, index]
56
57 conds = conds.transpose(1, 2) # [B, mel_dim, T_mel]
58
59 # 6. Flow mask preparation
60 # Flow expects shape [B, T]
61 flow_mask = (~make_pad_mask(h_lengths)).to(device)
62
63 # 7. Flow loss computation
64 loss, _ = self.decoder.compute_loss(
65     feat.transpose(1, 2).contiguous(),
66     flow_mask.unsqueeze(1),
67     h.transpose(1, 2).contiguous(),
68     embedding,
69     cond=conds,
70     streaming=streaming,
71 )
72
73 return {'loss': loss}

```

Listing 2: Modified

C Subjective Evaluation

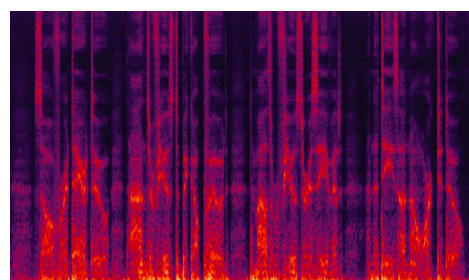
We developed a web-based automated subjective evaluation interface (Figure 3) to collect real-time listening ratings. Twenty audio stimuli (four per proposed method) were pre-loaded into four Latin-square sequences to minimise order and carry-over effects; the sequence was automatically selected according to the participant ID entered at the start of the session. Each trial started with the automatic playback of one stimulus (24 kHz, 16-bit WAV; approximately 12 s). Immediately after playback, a five-point Likert scale (0 = very bad, 5 = excellent) appeared on the screen and the participant had 5 s to tap the desired score on a touch device or click with a mouse. Responses were timestamped and written to a CSV file that was automatically downloaded at the end of the session.

D Case Study

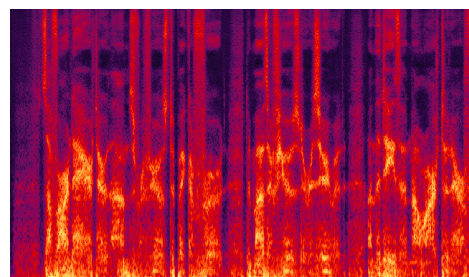
Sample spectrograms are provided above; additional examples and the full demo are available in the supplementary material.



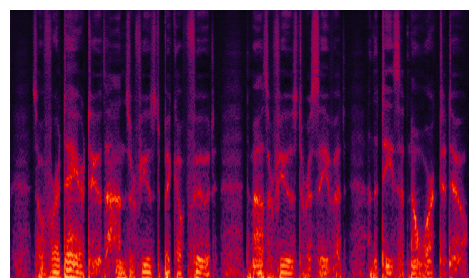
Figure 4: User interface of the listening test.



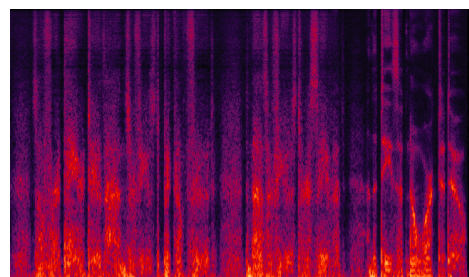
(a) CosyWhisper Output



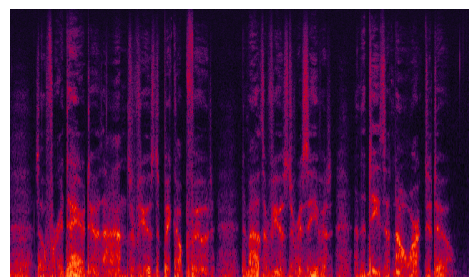
(b) CosyVoice3 Finetuned on WhispReal Output



(c) CosyVoice3 Repeat Output



(d) WhispSynth Output



(e) WhispReal Output

Figure 5: Sample spectrograms. (We invite reviewers to see the supplementary demo for more examples.)