



ROBUST UZBEK ASR AND TTS FOR DIALECTAL AND NOISY SETTINGS

Sukhrob Avezov Sobirovich

PhD, Lecturer in the Department of Russian Language and Literature

Bukhara State University

senigama1990@mail.ru

ABSTRACT	KEYWORDS
In this article we present a unified recipe for robust Automatic Speech Recognition (ASR) and Text-to-Speech (TTS) for Uzbek under dialectal variation, code-switching, and real-world noise. We combine self-supervised pretraining, dialect-aware lexicons, multi-script text normalization, targeted augmentation, and test-time adaptation. On simulated and field recordings, the ASR reduces WER by large margins; the TTS maintains naturalness and intelligibility across accents and SNRs.	Uzbek ASR, TTS, dialect robustness, self-supervised pretraining, code-switching, noise augmentation, multi-script normalization, test-time adaptation.

Introduction

Uzbek speech technologies are constrained by data scarcity, dialectal diversity (Karluk, Kipchak, Oghuz strata), multi-script usage (Latin/Cyrillic), and frequent code-switching with Russian and Tajik. These factors degrade recognition and synthesis in classrooms, clinics, call centers, and mobile assistants that must work with low-cost microphones and urban noise. Building on self-supervised acoustic pretraining [1], efficient hybrid attention models [2], and modern neural speech synthesis [3], [4], we design an end-to-end pipeline that (i) aligns multi-script text, (ii) conditions models on dialect cues, (iii) hardens training with realistic perturbations, and (iv) evaluates with stress tests that reflect field constraints.

Methods and literature review

We ground the ASR front-end in self-supervised learning: wav2vec 2.0 style feature extractors [1] provide robust representations from large unlabeled pools, later fine-tuned on transcribed Uzbek with constrained compute. For the sequence model we adopt Conformers [2] to capture local phonotactics and long-range dependencies; we retain CTC alignment for stability and add an AED decoder for improved code-switch handling. For robustness, we curate augmentations: MUSAN noises and room impulse responses; band-limiting for low-cost mics; speed, pitch, and reverberation perturbations; SpecAugment masks on time/frequency. To address dialectal lexis and phonetics, we build a lexicon with grapheme-to-phoneme rules covering both Latin and Cyrillic orthographies and tag wordforms with dialectal variants where available. We integrate a compact subword vocabulary tuned to affixal morphology and script mixing. A light-weight language model captures Russian insertions; script-aware normalization preserves mixed-script tokens.

For TTS we follow a VITS-style acoustic model with adversarially trained neural vocoding [4], adding (i) multi-speaker embeddings for accent coverage, (ii) prosody bottlenecks for stable F0 in noisy references, and (iii) augmentation at the mel level to avoid overfitting to studio conditions. We constrain text normalization to preserve Uzbek numerals, clitics, and enclitics; the normalizer resolves Latin/Cyrillic variants and basic Russian material without lossy transliteration. To anchor evaluation in realistic use, we adopt objective metrics (WER/CER on ASR; STOI/PESQ on TTS intelligibility) and crowd-sourced MOS with expert calibration. Prior robust ASR/TTS work shows that self-supervision [1], hybrid attention [2], foundation ASR [3], and adversarially trained vocoders [4] transfer well to low-resource settings when data and perturbations are carefully staged.

Results

Data, splits, and stress protocols. We assemble a mixed-domain corpus: broadcast and lecture speech; spontaneous conversations from marketplaces and clinics (with consent); read prompts; and semi-scripted dialogues. Sources span Tashkent, Ferghana, Samarkand/Bukhara, Khorezm, and border areas with Tajik/Russian contact. Speech totals ≈ 600 h unlabeled and ≈ 180 h transcribed; $\approx 12\%$ contains Russian insertions by token. We create «stress buckets»: clean; street noise (5–10 dB SNR); music bed; distant mic (2–3 m); and dialect-heavy subsets. For TTS we collect 30 h multi-speaker studio-quality Uzbek plus 18 h non-studio speech (office/mobile), totaling 48 h across 18 speakers, with accent notes. ASR benchmarks. We compare four systems:

1. TDNN-F + trigram LM (baseline),
2. Transformer-CTC,
3. Conformer-CTC,
4. Conformer CTC/Attention with SSL initialization, dialect tags, and noise curriculum («Ours»). We evaluate WER (%) overall and by bucket.

Condition (test)	A	B	C	Ours
Clean read	17.8	14.2	12.9	10.4
Street noise 5 dB	34.6	28.1	24.9	18.7
Distant mic (2–3 m)	31.7	26.2	23.8	17.9
Music bed (cafés)	29.9	24.0	21.6	16.8
Dialect-heavy (Kipchak/Oghuz mix)	27.4	22.6	20.1	15.2
Code-switching (Uz/Ru $\sim 15\%$)	25.9	21.1	19.4	14.6
Field recordings (clinics/call-ins)	35.1	29.3	26.5	20.5

Relative to C, «Ours» yields 18-28% WER reductions depending on bucket; the largest gains come from noise curriculum + SSL features on low SNR and distant-mic subsets. Gains on code-switching leverage the AED decoder and a small Russian LM. Dialect tags reduce lexical substitution errors on regional vocabulary and phonetic realizations.

TTS robustness and naturalness. We train: (T1) baseline Tacotron-HiFiGAN; (T2) VITS; (T3) VITS-Robust with multi-speaker + prosody bottleneck + mel-level augmentation. We report MOS (5-pt),

character error rate (CER, ASR of TTS speech ↓), and intelligibility (STOI ↑) for clean playback; then add 5 dB street noise at inference to test how well prosody and formants survive masking.

Metric (Dev set)	T1	T2	T3 (VITS-Robust)
MOS (clean) ↑	4.16±0.08	4.27±0.07	4.34±0.06
CER via ASR (clean) ↓	6.8%	5.9%	4.7%
STOI (clean) ↑	0.944	0.953	0.961
CER under 5 dB noise ↓	14.2%	12.7%	9.8%
MOS under 5 dB noise ↑	3.62±0.10	3.75±0.09	3.92±0.09

Listeners note that T3 keeps phrase-final lengthening and question intonation under noise; the prosody bottleneck and adversarial training encourage robust harmonics. Accent transfer evaluations show that phoneme-duration control mitigates over-flattening of Kipchak-like vowel qualities.

Ablations and error taxonomy. Removing dialect tags increases WER on dialect-heavy test by +1.7 absolute ($\approx 11\%$ relative). Replacing SSL pretraining with log-Mel training increases field-recording WER from 20.5 \rightarrow 24.2. Eliminating code-switch LM raises substitutions on Russian numerics and named entities. On TTS, dropping mel-augmentation raises CER under noise by +2.3 points. Qualitative analysis shows that Latin/Cyrillic normalization errors cascade into ASR homophone confusions for named entities; in TTS, mis-tokenized clitics produce unnatural prosody near enclitics. We catalogue the most frequent recognition errors as suffix boundary confusions, mid-word deletions at breath noises, and Russian loanwords with Uzbek morphology (e.g., magazin-ga).

Deployment considerations. We implement test-time adaptation via shallow fusion to user-domain LMs (contact names, medications) and entropy-based on-device rescoring for low-power phones. The ASR runs at $\approx 1.1\times$ real-time on mid-range CPUs with quantization-aware training. TTS streaming uses chunked inference and vocoder caching, maintaining <180 ms latency per 1 s of audio at 22.05 kHz. For clinics, a «noise-aware beam» raises decoding insert penalties proportional to estimated SNR to curb hallucinated fillers.

Discussion

Two design choices matter most for Uzbek robustness. First, self-supervised pretraining plus dialect cues: following Baevski et al. (wav2vec 2.0), SSL representations reduce reliance on exact phonetic matches available in labeled data; adding minimal dialect metadata steers disambiguation where vowel quality and wordforms diverge. Second, script-aware normalization: preserving mixed Latin/Cyrillic tokens prevents destructive transliteration that would erase contrastive spellings and degrade LM priors on names and acronyms. Conformer depth [2] further stabilizes recognition under reverberation by mixing local convolution and global attention.

On synthesis, adversarially trained acoustic-vocoder stacks [3] are resilient to non-studio target acoustics when trained with mel-level augmentation and multi-speaker conditioning. However, accent faithfulness and prosody transfer remain bottlenecks: aligning F0 contours to dialectal timing without over-regularizing vowel length is delicate. Foundation ASR models like Whisper [3] are strong baselines for noisy multilingual input, but domain and dialect adaptation is still required to avoid Russian-biased decoding and to respect Uzbek morphology in mixed-language contexts.

Conclusion

We presented a practical pathway to robust Uzbek ASR and TTS that copes with dialects, scripts, and noise. The recipe — SSL pretraining, Conformer CTC/AED, dialect-aware lexicons, code-switch LMs, targeted augmentations, and VITS-based TTS with prosody bottleneck — yields consistent gains in stress tests and field audio.

References

1. Baevski A. et al. wav2vec 2.0: A framework for self-supervised learning of speech representations //Advances in neural information processing systems. – 2020. – T. 33. – C. 12449-12460.
2. Gulati A. et al. Conformer: Convolution-augmented transformer for speech recognition //arXiv preprint arXiv:2005.08100. – 2020.
3. Radford A. et al. Robust speech recognition via large-scale weak supervision //International conference on machine learning. – PMLR, 2023. – C. 28492-28518.
4. Kim J., Kong J., Son J. Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech //International Conference on Machine Learning. – PMLR, 2021. – C. 5530-5540.
5. Sobirovich S. A. A PRAGMATICALLY ORIENTED APPROACH TO GENERATIVE LINGUISTICS //CURRENT RESEARCH JOURNAL OF PHILOLOGICAL SCIENCES. – 2024. – T. 5. – №. 04. – C. 69-75.