

# Do multilingual audio representation models demonstrate greater robustness to domain shift in code-switched sp

Assignee Research

June 10, 2026

## Abstract

Pre-trained models for automatic speech recognition (ASR) and speech enhancement (SE) have exhibited remarkable capabilities under matched noise and channel conditions. However, these models often suffer from severe performance degradation when confronted with domain shifts, particularly in the presence of unseen noise and channel distortions. In view of this, we in this paper present URSA-GAN, a unified and domain-aware generative framework specifically designed to mitigate mismatches in both noise and channel conditions. URSA-GAN leverages a dual-embedding architecture that consists of a noi

## 1 Introduction

This paper examines: Universal Robust Speech Adaptation for Cross-Domain Speech Recognition and Enhancement. Research question: Do multilingual audio representation models demonstrate greater robustness to domain shift in code-switched speech recognition tasks than monolingual counterparts when evaluated on standard speech benchmarks?.

## 2 Methodology

Systematic literature search across multiple databases yielded 15 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.8/10.

## 3 Results

15 papers retrieved. 10 claims extracted; 0 independently verified. Quality review score: 3.8/10.

## 4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

## 5 Extracted Claims

Claim	Verified	Confidence
URSA-GAN consistently demonstrates robustness and generalization in complex real-world scenarios.	×	0.10
PESQ estimates the perceptual quality of enhanced speech by comparing it to a clean reference, correlating well with human	×	0.04
STOI measures speech intelligibility by analyzing temporal and spectral similarity between the enhanced and reference signals	×	0.03
The mean opinion score (MOS) was adopted to assess simulated data realism.	×	0.01
The topline model, trained directly on labeled target-domain data, serves as an upper bound.	×	0.05
URSA-GAN shows a PESQ score of 3.16 and STOI of 95.30% on the VBD dataset.	×	0.09
DEMUCS shows a PESQ score of 3.05 and STOI of 95.2% on the VBD dataset.	×	0.01
Conventional domain adaptation approaches often rely on a significant amount of labeled target-domain data or entail complex	×	0.05
Data simulation has emerged as a viable alternative solution for domain adaptation in both ASR and SE.	×	0.05
Most current simulation techniques primarily focus on capturing broad domain properties and often neglect fine-grained, domain-specific	×	0.06

## References

- <http://arxiv.org/abs/2105.14779v2>
- <http://arxiv.org/abs/2602.04307v2>

- <http://arxiv.org/abs/2006.08870v1>