

# Does incorporating explicit phoneme alignment into UniSpeech-derived representations improve robustness in low

Assignee Research

June 10, 2026

## Abstract

Domain-specific languages that use a lot of specific terminology often fall into the category of low-resource languages. Collecting test datasets in a narrow domain is time-consuming and requires skilled human resources with domain knowledge and training for the annotation task. This study addresses the challenge of automated collecting test datasets to evaluate semantic search in low-resource domain-specific German language of the process industry. Our approach proposes an end-to-end annotation pipeline for automated query generation to the score reassessment of query-document pairs. To overc

## 1 Introduction

This paper examines: Automated Collection of Evaluation Dataset for Semantic Search in Low-Resource Domain Language. Research question: Does incorporating explicit phoneme alignment into UniSpeech-derived representations improve robustness in low-resource languages on LibriSpeech and CommonVoice benchmarks compared to baseline self-supervised models?.

## 2 Methodology

Systematic literature search across multiple databases yielded 11 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.5/10.

## 3 Results

11 papers retrieved. 0 claims extracted; 0 independently verified. Quality review score: 3.5/10.

## 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.

## References

- <http://arxiv.org/abs/2412.10008v1>
- <http://arxiv.org/abs/2501.05260v1>
- <http://arxiv.org/abs/2407.13292v2>