

# Multimodal Alignment Strategies in Self-Supervised Speech Representation Learning for Zero-Shot Transfer to Unseen Languages

Assignee Research

June 13, 2026

## Abstract

This paper proposes a zero-shot text-to-speech (TTS) conditioned by a self-supervised speech-representation model acquired through self-supervised learning (SSL). Conventional methods with embedding vectors from x-vector or global style tokens still have a gap in reproducing the speaker characteristics of unseen speakers. A novel point of the proposed method is the direct use of the SSL model to obtain embedding vectors from speech representations trained with a large amount of data. We also introduce the separate conditioning of acoustic features and a phoneme duration predictor to obtain the

## 1 Introduction

This paper examines: Zero-shot text-to-speech synthesis conditioned using self-supervised speech representation model. Research question: How do multimodal alignment strategies in self-supervised speech representation learning affect zero-shot transfer performance on unseen languages?.

## 2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.9/10.

## 3 Results

13 papers retrieved. 10 claims extracted; 8 independently verified. Quality review score: 7.9/10.

## 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
The proposed method uses a self-supervised learning (SSL) model to obtain embedding vectors instead of x-vectors.	✓	0.22
GenerSpeech requires fine-tuning the SSL model for speaker and emotion recognition.	✓	0.17
Under the parallel condition, the proposed method achieved lower Mean Absolute Error (MAE) for log melspectrograms than	✓	0.18
Under the parallel condition, the proposed method achieved lower Root Mean Square Error (RMSE) for phoneme duration than	✓	0.22
Separate conditioning with LSTM-based aggregation yielded the best performance for both w2v2 and HuBERT models compared	✓	0.16
Using separate conditioning with LSTM-based aggregation drastically improved the RMSEs of phoneme duration.	✓	0.22
For the HuBERT LARGE model trained on LibriLight under parallel conditions, the log mel-spectrogram MAE was 1.12.	×	0.09
For the HuBERT LARGE model trained on LibriLight under parallel conditions, the phoneme duration RMSE was 16.1.	×	0.08
Objective metrics under the non-parallel condition were calculated by comparing the target speaker's speech with the gen	✓	0.26
The dimension of the speech representation vector is defined as $D_{sr} = D_l \times L$ .	✓	0.18

## References

- <http://arxiv.org/abs/2007.04134v1>

- <http://arxiv.org/abs/2503.06273v2>
- <http://arxiv.org/abs/2304.11976v1>