

Comparative Analysis of Joint Audiovisual Self-Supervised and Supervised Speech Representations for Speaker Recognition and

Assignee Research

June 12, 2026

Abstract

The intuitive interaction between the audio and visual modalities is valuable for cross-modal self-supervised learning. This concept has been demonstrated for generic audiovisual tasks like video action recognition and acoustic scene classification. However, self-supervision remains under-explored for audiovisual speech. We propose a method to learn self-supervised speech representations from the raw audio waveform. We train a raw audio encoder by combining audio-only self-supervision (by predicting informative audio attributes) with visual self-supervision (by generating talking faces from au

1 Introduction

This paper examines: Learning Speech Representations from Raw Audio by Joint Audiovisual Self-Supervision. Research question: How does the performance of self-supervised speech representations learned from raw audio using joint audiovisual self-supervision compare to supervised representations on downstream tasks like speaker recognition and emotion detection, as measured by accuracy and generalization across unseen datasets?.

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 7.2/10.

3 Results

15 papers retrieved. 22 claims extracted; 17 independently verified. Quality review score: 7.2/10.

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 LRW dataset contains 500 different isolated words primarily from BBC recordings.	✓	0.23
Filtering the LRW dataset for videos with yaw, pitch, and roll restricted to a maximum of 10 degrees results in approxim	✓	0.24
The Speech Commands (SPC) dataset contains 64,727 total utterances of 30 different words by 1,881 speakers.	✓	0.26
The study compares proposed methods against PASE, APC, wav2vec, L1, and L1 + Odd baselines using code and pretrained mod	×	0.15
The L1 baseline method is based on log mel spectrograms rather than raw audio.	×	0.14
The first supervised baseline uses 39-dimensional MFCCs consisting of 13 coefficients, 13 deltas, and 13 delta-deltas.	✓	0.16
The second supervised baseline is a fully supervised 1D Resnet18 model trained from scratch directly on the evaluation d	✓	0.22
Downstream evaluation is performed on isolated word classification using the Speech Commands (SPC) and Lip Reading in th	✓	0.15
The downstream classifier architecture consists of a 2-layer BiGRU with 256 units per layer followed by a linear layer.	✓	0.15
The number of target classes for the downstream classifier is 30 for SPC and 500 for LRW.	×	0.15
Models are finetuned for downstream classification for 50 epochs.	×	0.08
The learning rate schedule is 0.0001 for the first 40 epochs and 0.00001 for the last 10 epochs.	✓	0.15
Training uses standard softmax plus cross entropy loss.	×	0.09
The proposed audio encoder architecture is a 1D Resnet18.	✓	0.16
The audio encoder takes a 16 kHz raw audio waveform as input.	✓	0.16
The audio encoder outputs a 512-dimensional audio feature vector for every timestep.	✓	0.18
The output sample rate of the audio encoder is 25 audio feature vectors per second.	✓	0.23
The audio encoder output rate of 25 vectors per second matches the 25 FPS video 4rate in the LRW dataset to enable one-to	✓	0.21
Contemporary self-supervised methods such as Alwassel et al. (2019) and Patrick et al. (2020) use a 2D Resnet18 audio en	✓	0.33
The visual self-supervision model generates a talking lip video from a still image and corresponding audio	✓	0.17

References

- <http://arxiv.org/abs/2208.05445v1>
- <http://arxiv.org/abs/2402.06923v1>
- <http://arxiv.org/abs/2007.04134v1>