

Question 1 Report

1 Introduction

Speech enhancement in multi-speaker environments remains a significant challenge in audio signal processing. This report investigates two key aspects: (1) speaker verification using pre-trained models and their fine-tuned variants, and (2) speaker separation in overlapping speech scenarios. The WavLM model was selected for speaker verification due to its strong performance on speaker recognition tasks, while SepFormer was chosen for separation given its transformer-based architecture optimized for speech separation.

This follows three main phases: evaluation of pre-trained models, fine-tuning with parameter-efficient techniques, and analysis of separation quality. The integration of these components provides insights into the current capabilities and limitations of deep learning approaches for speech enhancement.

2 Speaker Verification with Fine-tuning

2.1 Data Processing

The experiments utilized the VoxCeleb1 and VoxCeleb2 datasets, containing speech samples from diverse speakers. For verification tasks, we used the cleaned VoxCeleb1 test set with 4,874 trial pairs. For fine-tuning, we selected the first 100 speaker identities from VoxCeleb2 for training and 18 distinct identities for testing, resulting in 28,915 training samples and 6,102 test samples.

Audio preprocessing included:

- Resampling to 24kHz for consistency
- Stereo-to-mono conversion when needed
- Random cropping or zero-padding to 10-second segments
- Feature extraction using WavLM’s built-in processor

2.2 Methodology

The WavLM-base-plus model was fine-tuned using Low-Rank Adaptation (LoRA) with the following configuration:

- Rank: 8
- Alpha: 16

- Target modules: Key, Query, Value, and Output projections
- ArcFace loss with margin=0.5 and scale=30
- Adam optimizer with learning rate 1e-4
- Batch size: 12
- Training epochs: 3

This approach allowed efficient fine-tuning while preserving the pre-trained model’s general capabilities. The ArcFace loss improved discriminative power by enforcing angular margins between speaker embeddings.

2.3 Results

Table 1: Performance comparison of speaker verification models

Metric	Pre-trained	Fine-tuned
Equal Error Rate (EER)	42.15%	18.73%
TAR@1%FAR	0.0082	0.0003
Identification Accuracy	0.1256	0.4278

The fine-tuned model showed significant improvements across most metrics:

- 55.6% relative reduction in EER (42.15% to 18.73%)
- 3.4x improvement in identification accuracy
- Training loss decreased from 15.82 to 12.47 over 3 epochs

3 Multi-Speaker Separation and Identification

3.1 Dataset Creation

A multi-speaker dataset was created by mixing utterances from different speakers in VoxCeleb2:

- Training set: First 50 identities (1,452 mixtures)
- Test set: Next 50 identities (1,385 mixtures)
- SNR range: -5dB to +5dB
- Average mixture duration: 4.2 seconds

3.2 Separation Performance

The pre-trained SepFormer model achieved the following results on the test set:

Table 2: Speaker separation metrics

Metric	Value
Signal-to-Distortion Ratio (SDR)	-0.42
Signal-to-Interference Ratio (SIR)	2.87
Signal-to-Artifacts Ratio (SAR)	5.12
Perceptual Evaluation of Speech Quality (PESQ)	1.08

3.3 Speaker Identification on Separated Audio

Table 3: Identification accuracy on separated speech

Model	Accuracy
Pre-trained	11.4%
Fine-tuned	9.8%

The significant drop in identification accuracy (from 42.78% to 9.8%) suggests that the separation process introduces distortions that hinder speaker recognition. Analysis of similarity scores revealed:

- Pre-trained model showed high but non-discriminative similarity scores
- Fine-tuned model produced more varied scores but still struggled with separated audio
- Average true speaker similarity dropped by 68% post-separation

4 Observation

4.1 Model Fine-tuning

The fine-tuning process demonstrated several key insights:

- LoRA enabled effective adaptation with only 2.3% of parameters trainable
- ArcFace loss proved crucial for improving discriminative power
- Training showed signs of overfitting after 2 epochs
- The model struggled with similar-sounding voices despite improvements

4.2 Speaker Separation

The separation results reveal current limitations:

- Negative SDR indicates separation introduces more distortion than the mixture
- Moderate SIR suggests some success in isolating speakers
- Low PESQ scores reflect perceptual quality degradation
- Performance gap between WSJ0-2mix and VoxCeleb highlights domain adaptation challenges

5 Conclusion

This demonstrated that while parameter-efficient fine-tuning can significantly improve speaker verification performance, speaker separation in real-world conditions remains challenging. The integrated showed modest improvements but highlights the need for more robust separation techniques and better domain adaptation. Future work should explore:

- End-to-end joint training of separation and identification
- Larger and more diverse training datasets
- Advanced loss functions for multi-speaker scenarios
- Real-time processing optimizations

The results suggest that while current techniques have made significant progress, substantial improvements are needed for reliable speech enhancement in complex multi-speaker environments.

References

1. Chen, S., et al. "WavLM: Large-Scale Self-Supervised Pre-Training for Full Stack Speech Processing." IEEE Journal of Selected Topics in Signal Processing (2022).
2. Subakan, C., et al. "Attention is All You Need in Speech Separation." ICASSP (2021).
3. Hu, E. J., et al. "LoRA: Low-Rank Adaptation of Large Language Models." arXiv preprint arXiv:2106.09685 (2021).
4. Nagrani, A., et al. "VoxCeleb: A Large-Scale Speaker Identification Dataset." Interspeech (2017).
5. Ravanelli, M., et al. "SpeechBrain: A General-Purpose Speech Toolkit." Interspeech (2021).