

Multi-Speaker Speech Enhancement and Separation

Question 1

Harishankar Nagar (D24CSA003)

Abstract

In this report, a comprehensive analysis of multi-speaker speech enhancement using speaker verification models and separation techniques is presented. The effectiveness of pre-trained and fine-tuned models for speaker verification tasks was evaluated using the VoxCeleb dataset. Speaker separation in multi-speaker environments was explored using the SepFormer model. The performance was assessed using metrics such as Equal Error Rate (EER), Signal to Interference Ratio (SIR), Signal to Artifacts Ratio (SAR), Signal to Distortion Ratio (SDR), and Perceptual Evaluation of Speech Quality (PESQ). The findings highlight the challenges and potential solutions for enhancing speech in multi-speaker scenarios.

1 Introduction

Speech enhancement and speaker separation represent critical challenges in audio processing, particularly in environments with multiple speakers. This assignment focuses on enhancing speech intelligibility and quality in multi-speaker scenarios using state-of-the-art deep learning techniques. Pre-trained models were leveraged and fine-tuned for specific tasks, with their performance evaluated using various metrics.

The assignment consists of multiple tasks: (1) speaker verification using pre-trained models, (2) fine-tuning these models with Low-Rank Adaptation (LoRA) [1] and ArcFace loss [2], (3) creating multi-speaker datasets and performing speaker separation, and (4) evaluating the quality of separated speech.

2 Task 2: Speaker Verification with Fine-tuning

2.1 Data Processing

For this task, the VoxCeleb1 [3] and VoxCeleb2 [4] datasets were utilized, which contain speech recordings from a diverse range of speakers. The WavLM base plus model [5] was selected for speaker verification and the following data processing steps were performed:

- Trial pairs for verification were loaded from the VoxCeleb1 (cleaned) dataset.
- 100 identities (sorted in ascending order) from VoxCeleb2 were selected for training and 18 identities for testing.
- Audio files were processed to ensure consistent sampling rate (16kHz) and stereo to mono conversion was handled when needed.
- Data loaders were created for efficient batch processing during training and evaluation.

The VoxCeleb2 dataset was used for fine-tuning the model, with the first 100 identities allocated for training and the remaining 18 identities for testing. This resulted in 29,831 training samples and 6,406 testing samples.

2.2 Methodology

A two-stage approach was implemented:

1. **Pre-trained Model Evaluation:** First, the performance of the pre-trained WavLM base plus model was evaluated on speaker verification tasks.
2. **Fine-tuning with LoRA and ArcFace Loss:** The model was then fine-tuned using Low-Rank Adaptation (LoRA) and ArcFace loss, a technique that enhances the discriminative power of deep embeddings by introducing an angular margin.

For the fine-tuning process, LoRA was applied with a rank of 8 and a scaling factor of 16, targeting the attention modules (key, query, value, and output projections). The ArcFace loss was implemented with a scale of 30.0 and a margin of 0.5. The Adam optimizer [6] was used with a learning rate of 1e-4, and the model was trained for 3 epochs with a batch size of 12.

Implementation was done using PyTorch [7] with the Hugging Face Transformers library [8] for accessing the pre-trained models and the PEFT library [9] for implementing LoRA fine-tuning. The torchaudio library [10] was used for audio processing and the scikit-learn library [11] for evaluation metrics.

2.3 Results

The performance comparison between the pre-trained and fine-tuned models is summarized below:

Metric	Pre-trained Model	Fine-tuned Model
Equal Error Rate (EER)	47.94%	19.25%
True Acceptance Rate @ 1% FAR	0.0060	0.0001
Speaker Identification Accuracy	0.1099	0.4364

Table 1: Performance comparison of pre-trained and fine-tuned models

A significant improvement was observed in the Equal Error Rate (EER) after fine-tuning, with a reduction from 47.94% to 19.25%. This represents a substantial enhancement in the model’s ability to verify speakers correctly. Similarly, the speaker identification accuracy increased substantially from 0.1099 to 0.4364, indicating that the fine-tuned model was much more effective at correctly identifying speakers. However, the True Acceptance Rate at 1% False Acceptance Rate (TAR@1%FAR) decreased slightly from 0.0060 to 0.0001, suggesting a trade-off in the model’s performance characteristics.

These results demonstrate that the fine-tuning approach with LoRA and ArcFace loss was effective in improving the overall speaker verification and identification capabilities of the model. The reduction in EER is particularly noteworthy, as it indicates a better balance between false acceptances and false rejections. The improvement in identification accuracy further supports the effectiveness of the fine-tuning process.

3 Task 3: Multi-Speaker Separation and Identification

3.1 Data Creation and Processing

For this task, a multi-speaker scenario dataset was created by mixing utterances from different speakers in the VoxCeleb2 dataset. The following steps were taken:

- The first 50 identities (sorted in ascending order) from VoxCeleb2 were used to create a multi-speaker training scenario.
- The next 50 identities were used to create a multi-speaker testing scenario.
- Audio files from pairs of speakers were mixed with varying Signal-to-Noise Ratios (SNR) to create realistic overlapping speech.
- For each pair, the original clean sources were preserved for evaluation purposes.

In total, 50 mixed samples were created for both the training and testing sets. Each mixed sample contained speech from two different speakers, with the mixing process carefully controlled to ensure a balanced representation of both speakers. The approach for mixing was inspired by the methodology used in the LibriMix dataset [12].

3.2 Methodology

The SepFormer model [13] was used for speaker separation, which is a state-of-the-art approach based on the Transformer architecture. The pre-trained SepFormer model was applied to the created test set, and the separated speech was evaluated using multiple metrics. Additionally, both the pre-trained and fine-tuned speaker identification models from Task 2 were used to identify which enhanced speech corresponds to which speaker.

The process involved the following steps:

1. The SepFormer model was loaded from the pre-trained checkpoint (speechbrain/sepformer-wsj02mix) using the SpeechBrain library [14].
2. Each mixed audio sample was processed through the model to separate the overlapping speakers.
3. The resulting separated sources were saved and evaluated against the original clean sources.
4. Speaker identification was performed on the separated sources using both the pre-trained and fine-tuned models from Task 2.

For evaluation, the mir_eval library [15] was used to compute the BSS (Blind Source Separation) metrics including SDR, SIR, and SAR. The PESQ implementation from the pesq library was used to evaluate the perceptual quality of the separated audio.

3.3 Results

The performance of the SepFormer model for speaker separation is summarized below:

For speaker identification on the separated audio, the following results were obtained:

The SepFormer model exhibited mixed performance in separating speakers. The negative SDR value (-0.1555) indicates that some distortion was introduced in the separation process, potentially

Metric	Value
Signal to Distortion Ratio (SDR)	-0.1555
Signal to Interference Ratio (SIR)	3.0043
Signal to Artifacts Ratio (SAR)	5.3592
Perceptual Evaluation of Speech Quality (PESQ)	1.1399

Table 2: Performance metrics for speaker separation using SepFormer

Model	Rank-1 Identification Accuracy
Pre-trained	12.00%
Fine-tuned	10.00%

Table 3: Speaker identification accuracy on separated audio

making the separated sources less recognizable than the original mixture. However, the positive SIR (3.0043) suggests that the model was somewhat effective at reducing interference between speakers. The SAR value (5.3592) indicates moderate performance in reducing artifacts in the separated audio. The PESQ score (1.1399) is relatively low, suggesting that the perceptual quality of the separated speech was significantly degraded compared to the original clean speech.

For speaker identification, both models performed poorly on the separated audio, with accuracies of 12.00% and 10.00% for the pre-trained and fine-tuned models, respectively. This significant drop in performance compared to Task 2 (where the fine-tuned model achieved 43.64% accuracy) suggests that the separated audio contained substantial distortions that made speaker identification challenging.

Analyzing the similarity scores to diagnose the low identification accuracy revealed:

- Pre-trained model average true speaker similarity: 0.9388
- Pre-trained model average predicted speaker similarity: 0.9570
- Fine-tuned model average true speaker similarity: 0.1888
- Fine-tuned model average predicted speaker similarity: 0.3856

These numbers indicate that while the pre-trained model produced high similarity scores overall (suggesting overconfidence), the fine-tuned model had much lower scores, with the predicted speaker still having a higher score than the true speaker on average.

4 Discussion

4.1 Task 2: Model Fine-tuning

The significant improvement in EER and identification accuracy after fine-tuning demonstrates the effectiveness of the LoRA technique combined with ArcFace loss. The reduction in EER from 47.94% to 19.25% represents a 59.84% relative improvement, indicating that the fine-tuned model is much better at distinguishing between same-speaker and different-speaker pairs.

The model showed the following training progression:

- Epoch 1: Train Loss = 16.4628, Val Loss = 22.2481, Accuracy = 0.0072

- Epoch 2: Train Loss = 13.7926, Val Loss = 23.2158, Accuracy = 0.0184
- Epoch 3: Train Loss = 12.5018, Val Loss = 23.3379, Accuracy = 0.0089

The decreasing training loss indicates that the model was improving on the training data, while the increasing validation loss suggests some overfitting. However, the final identification accuracy of 0.4364 on the test set demonstrates that the model still generalized well despite this potential overfitting.

4.2 Task 3: Speaker Separation

The speaker separation results highlight the challenges of separating overlapping speech, particularly in non-ideal conditions. The negative SDR value suggests that the separated signals contain more distortion than the original mixture, which could be due to the complexity of the overlapping speech patterns in the VoxCeleb2 dataset compared to the WSJ0-2mix dataset [16] that the SepFormer was originally trained on.

The significant drop in speaker identification accuracy on the separated audio (from 43.64% to 10.00% for the fine-tuned model) underscores the difficulty of the task. The separation process introduces distortions that make it harder for the identification model to recognize speakers accurately. This is further supported by the similarity score analysis, which shows that the fine-tuned model had much lower confidence in its predictions overall, with the predicted speaker still having a higher similarity score than the true speaker.

5 Conclusion

In this assignment, the effectiveness of pre-trained and fine-tuned models for speaker verification and identification in multi-speaker environments was evaluated. The fine-tuning approach with LoRA and ArcFace loss significantly improved the speaker verification performance, as evidenced by the substantial reduction in EER and the increase in identification accuracy.

However, the performance of these models on separated speech was considerably lower, highlighting the challenges of speaker identification in multi-speaker scenarios. The SepFormer model showed limited effectiveness in separating overlapping speech from the VoxCeleb2 dataset, as indicated by the negative SDR and low PESQ scores.

These findings suggest that while significant progress has been made in speaker verification and separation techniques, there is still room for improvement, particularly in handling real-world multi-speaker scenarios. Future work could focus on developing more robust separation models specifically trained on diverse multi-speaker datasets like VoxCeleb and exploring advanced techniques for speaker identification in noisy and overlapping conditions.

References

- [1] E. J. Hu, Y. Shen, P. Wallis, *et al.*, “Lora: Low-rank adaptation of large language models,” *arXiv preprint arXiv:2106.09685*, 2021.
- [2] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, “Arcface: Additive angular margin loss for deep face recognition,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4690–4699.
- [3] A. Nagrani, J. S. Chung, and A. Zisserman, “Voxceleb: A large-scale speaker identification dataset,” in *Interspeech*, 2017, pp. 2616–2620.

- [4] J. S. Chung, A. Nagrani, and A. Zisserman, “Voxceleb2: Deep speaker recognition,” in *Interspeech*, 2018, pp. 1086–1090.
- [5] S. Chen, C. Wang, Z. Chen, *et al.*, “Wavlm: Large-scale self-supervised pre-training for full stack speech processing,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, no. 6, pp. 1505–1518, 2022.
- [6] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *International Conference on Learning Representations (ICLR)*, 2014.
- [7] A. Paszke, S. Gross, F. Massa, *et al.*, “Pytorch: An imperative style, high-performance deep learning library,” in *Advances in Neural Information Processing Systems*, 2019, pp. 8026–8037.
- [8] T. Wolf, L. Debut, V. Sanh, *et al.*, “Transformers: State-of-the-art natural language processing,” *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 38–45, 2020.
- [9] Hugging Face, *Peft: Parameter-efficient fine-tuning of billion-scale models on low-resource hardware*, <https://github.com/huggingface/peft>, 2023.
- [10] Y.-Y. Yang, M. Ravanelli, Y. Braun, S. K. Chung, K.-C. Peng, and M. Chang, “Torchaudio: Building blocks for audio and speech processing,” *arXiv preprint arXiv:2110.15018*, 2022.
- [11] F. Pedregosa, G. Varoquaux, A. Gramfort, *et al.*, “Scikit-learn: Machine learning in python,” *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [12] J. Cosentino, M. Pariente, S. Cornell, A. Deleforge, and E. Vincent, “Librimix: An open-source dataset for generalizable speech separation,” in *Interspeech*, 2020, pp. 2793–2797.
- [13] C. Subakan, M. Ravanelli, S. Cornell, M. Bronzi, and J. Zhong, “Attention is all you need in speech separation,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 21–25.
- [14] M. Ravanelli, T. Parcollet, P. Plantinga, *et al.*, “Speechbrain: A general-purpose speech toolkit,” in *Interspeech*, 2021.
- [15] C. Raffel, B. McFee, E. J. Humphrey, *et al.*, “Mir_eval: A transparent implementation of common mir metrics,” *Proceedings of the 15th International Society for Music Information Retrieval Conference*, 2014.
- [16] J. R. Hershey, Z. Chen, J. Le Roux, and S. Watanabe, “Deep clustering: Discriminative embeddings for segmentation and separation,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2016, pp. 31–35.