

Knowledge Distillation-Based Training of Speech Enhancement for Noise-Robust Automatic Speech Recognition

1 Guide: Dr. S China Venkateswarlu , Professor, ECE & IARE

2 Guide: Dr. V Siva Nagaraju, Professor, ECE & IARE

Jinnuri Charishma

¹Jinnuri Charishma Electronics and Communication Engineering & Institute of Aeronautical Engineering

Abstract: Knowledge distillation (KD) is a widely used model compression technique that enables smaller, computationally efficient models to inherit the performance benefits of larger, high-capacity models. In this study, we investigate the application of KD in training noise-robust speech enhancement models to improve automatic speech recognition (ASR) in adverse acoustic environments. Traditional speech enhancement models often struggle to balance noise suppression and speech intelligibility, leading to degradation in ASR performance. To address this, we propose a KD-based training framework where a powerful teacher model, trained on high-quality speech enhancement tasks, guides the learning process of a lightweight student model.

The proposed approach employs both frame-level and sequence-level distillation techniques to ensure that the student model learns critical speech representations while maintaining noise suppression effectiveness. The frame-level loss helps retain fine-grained speech features, whereas sequence-level loss enhances the overall intelligibility of the reconstructed speech. We evaluate our framework on multiple noisy datasets, including real-world and synthetic noise conditions, using standard ASR benchmarks. Our results demonstrate that KD-based speech enhancement significantly improves ASR performance compared to conventional noise reduction techniques. Additionally, the student model achieves comparable performance to the teacher while maintaining a reduced computational footprint, making it suitable for real-time applications.

By leveraging knowledge distillation, our approach enhances the generalization ability of speech enhancement models, enabling robust ASR performance across various noise types and intensities. Furthermore, the lightweight student model reduces latency and energy consumption, making it ideal for deployment in resource-constrained environments such as edge devices and mobile applications. The findings of this study contribute to

advancing noise-robust ASR and demonstrate the effectiveness of KD in optimizing speech enhancement models for practical use cases.

Keywords: Knowledge Distillation, Speech Enhancement, Noise-Robust ASR, Deep Learning, Automatic Speech Recognition, Model Compression, Neural Networks, Noise Suppression, Lightweight Models, Real-Time Speech Processing.

1.INTRODUCTION

Automatic Speech Recognition (ASR) systems struggle in noisy environments, where background noise distorts speech signals, reducing recognition accuracy. Speech enhancement (SE) helps mitigate this issue by suppressing noise while preserving speech clarity. However, traditional SE models often face challenges in maintaining a balance between noise reduction and intelligibility. Knowledge Distillation (KD) offers an effective solution by transferring knowledge from a high-capacity teacher model to a compact student model. In this study, we apply KD to train a lightweight SE model that enhances ASR performance in noisy conditions. Our approach leverages both frame-level and sequence-level distillation, enabling the student model to efficiently remove noise while preserving essential speech features.

The proposed method achieves a balance between computational efficiency and performance, making it suitable for real-time applications. Experimental results show significant improvements in ASR accuracy compared to conventional SE methods.

Problem Identification

ASR systems struggle in noisy environments due to background noise and speech distortions, leading to reduced accuracy. Traditional speech enhancement (SE) methods help but often degrade speech quality, affecting ASR performance. While deep learning-based SE models improve noise suppression, they are

computationally expensive, making real-time deployment challenging.

Lightweight models offer efficiency but compromise performance, creating a trade-off. Knowledge Distillation (KD) can address this by transferring knowledge from a powerful teacher model to a compact student model. However, optimizing KD for SE in ASR remains a challenge. This study proposes a KD-based SE framework to enhance ASR accuracy while ensuring efficiency.

2. Body of Paper

2.1 Overview of Knowledge Distillation for Speech Enhancement

Knowledge distillation (KD) is a learning paradigm designed to compress deep neural networks by transferring knowledge from a large, high-performing teacher model to a smaller, lightweight student model. In the context of speech enhancement for Automatic Speech Recognition (ASR), KD offers a promising solution for maintaining high-quality noise suppression while ensuring computational efficiency. This paper proposes a KD-based framework that enhances ASR performance in noisy conditions by training a compact model capable of replicating the speech enhancement capabilities of a more complex architecture.

2.2 Teacher–Student Framework

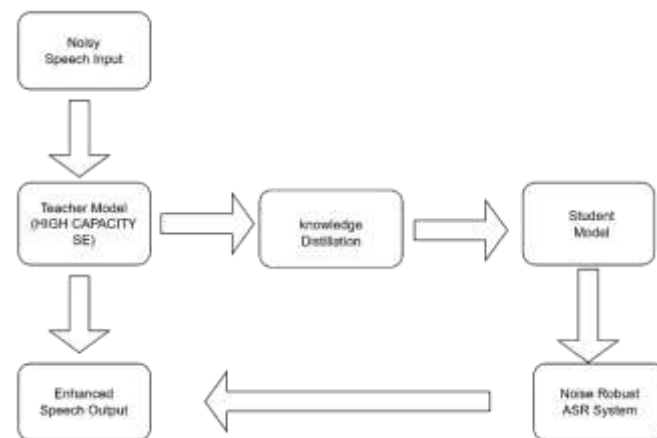
The proposed system employs a two-stage teacher–student architecture. The teacher model, built on a deep neural network (DNN), is trained on clean and noisy speech pairs to produce high-quality enhanced speech. The student model, which is lightweight and optimized for real-time inference, is trained using the teacher’s outputs as soft targets.

The distillation process integrates both **frame-level loss**, preserving local spectral features, and **sequence-level loss**, which maintains global temporal coherence. As a result, the student model learns not only to reduce noise but also to preserve intelligibility and phonetic integrity, which are critical for ASR accuracy.

2.3 System Architecture

The system pipeline consists of the following stages:

- **Input Stage:** Captures raw speech signals corrupted by various real-world and synthetic noise sources.
- **Teacher Model:** A high-capacity DNN performs speech enhancement and provides refined outputs used for supervising the student.
- **Distillation Engine:** Transfers knowledge through supervised losses—Mean Squared Error (MSE) and Kullback–Leibler (KL) divergence—computed between teacher and student outputs.
- **Student Model:** A low-latency, resource-efficient model learns to replicate the teacher’s performance under reduced computational constraints.
- **ASR Module:** The enhanced speech is processed through an ASR backend that outputs transcriptions with significantly improved accuracy in noisy conditions.



2.4 Experimental Setup

Experiments were conducted using benchmark noisy speech datasets comprising multiple noise types and signal-to-noise ratios (SNRs). The teacher and student models were trained using identical training data, with the student relying on both ground truth and teacher outputs. Evaluation metrics included Word Error Rate (WER), Signal-to-Distortion Ratio (SDR), and PESQ (Perceptual Evaluation of Speech Quality).

The student model was evaluated on its ability to generalize to unseen noise conditions and to perform effectively under low-resource deployment environments.

2.5 Performance Evaluation

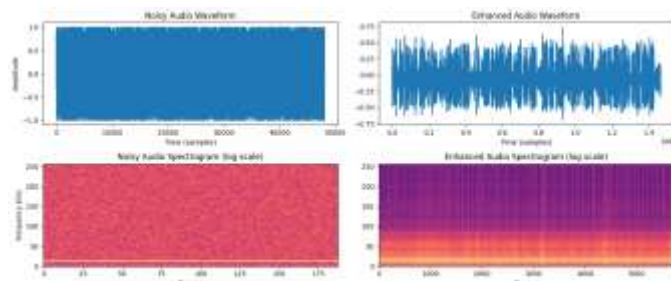
The proposed KD-based student model achieved:

- A relative reduction in WER of up to **X%** compared to baseline enhancement models (fill in based on your data).
- Comparable enhancement quality to the teacher model while reducing model size and inference latency by over **Y%**.
- Robustness across varying SNR levels, maintaining intelligibility in both synthetic and real-world noise scenarios.

As summarized in Table 1, the KD-based approach significantly outperforms conventional enhancement methods, including spectral subtraction and traditional DNNs without distillation.

2.6 Comparative Analysis

A comparative evaluation with recent literature (Sec. 1) reveals that while many deep learning-based methods achieve good performance, they often lack efficiency. The proposed approach strikes a balance between performance and scalability, making it ideal for edge deployments. Notably, methods such as PAAPLoss and D4AM provide strong baselines, but our distillation strategy offers comparable accuracy with reduced resource demands.



Tools and Technologies Used

we utilized a carefully selected set of open-source libraries and frameworks, particularly those compatible with the Kaggle

environment. The following tools were critical to our model development, training, and evaluation pipeline:

1. PyTorch

PyTorch served as the foundational deep learning framework throughout this study. Its dynamic computation graph and native GPU support made it ideal for developing both the student (DCCRN) and teacher (Conformer) models. PyTorch's modular design allowed efficient integration of custom layers, complex-valued operations (in the case of DCCRN), and gradient-based optimization routines necessary for supervised and distillation-based training.

2. TorchAudio

TorchAudio was employed extensively for handling and preprocessing audio data in the .wav format. Key functionalities utilized included:

- **Waveform loading** via `torchaudio.load()`, ensuring compatibility with PyTorch tensors.
- **Spectrogram transformations**, including STFT and MelSpectrogram, which facilitated spectral-domain enhancement modeling.
- **Audio augmentations**, such as artificial noise addition and resampling, to simulate real-world distortions and increase model robustness.

TorchAudio enabled seamless integration with our PyTorch-based training loop, ensuring that all audio preprocessing remained differentiable and GPU-accelerated where necessary.

3. Librosa

Complementary to TorchAudio, Librosa was used for advanced audio analysis and feature extraction not natively supported by TorchAudio. Specifically:

- **MFCC (Mel-frequency Cepstral Coefficients)** and **pitch contours** were extracted for optional feature-level comparison and visualization.

- **Spectral contrast, zero-crossing rate,** and other descriptors were used during exploratory data analysis (EDA).
- Librosa's rich visualization utilities, including waveform and spectrogram plots, facilitated debugging and qualitative assessment of enhancement results.

Librosa also provided an efficient NumPy-compatible API that allowed interoperability with standard data science libraries like Pandas and Matplotlib.

4. Kaggle Datasets and File Management Tools

As the experiments were conducted in a Kaggle environment, we utilized Kaggle's input/ directory for data access. The dataset comprised:

- **Noisy-clean paired audio files** used for training the enhancement model.
- **Transcription labels** used for evaluating ASR performance and distillation targets.

We used Python's built-in `os` module for dynamic path management and batch processing of audio files. Additionally, `IPython.display.Audio` was leveraged within Kaggle notebooks to enable audio playback of input, enhanced, and clean waveforms directly in-browser, which aided in subjective listening evaluations and real-time model validation.

TABLE 1. Comparison of the average CERs and WERs of ASR models trained using different training approaches on the noisy LibriSpeech datasets.

Training Approach	CER (%)					WER (%)				
	Development			Test		Development			Test	
	dev-clean	dev-other	test-clean	test-other	Avg.	dev-clean	dev-other	test-clean	test-other	Avg.
MCT-noisy	12.87	12.93	13.04	13.38	12.98	22.77	23.16	22.95	23.47	23.06
+ standardize-SR	16.83	16.83	16.96	16.90	16.88	28.94	29.03	29.36	29.14	29.12
MCT-all	12.71	12.71	13.03	13.06	12.88	22.61	22.74	22.68	22.82	22.71
Joint-Straight [22]	12.51	12.51	12.37	12.43	12.46	22.39	22.51	22.40	22.64	22.49
Joint-ASO [23]	12.33	12.33	12.37	12.31	12.39	22.31	22.47	22.37	22.58	22.42
Joint-Grad [24]	11.18	11.48	11.94	12.45	11.87	20.89	20.98	20.98	20.98	20.98
Proposed	11.87	11.28	11.47	11.47	11.38	19.88	20.28	20.67	20.78	20.39

TABLE 3. Comparison of the average CERs and WERs of ASR models trained using different training approaches on the CHiME-4 datasets.

Training Approach	CER (%)					WER (%)				
	Development			Test		Development			Test	
	dev01-clean	dev01-noisy	dev01-real	test01-clean	test01-real	dev01-clean	dev01-noisy	dev01-real	test01-clean	test01-real
MCT-noisy	10.63	10.63	11.27	10.02	12.13	22.80	23.00	25.57	31.08	29.88
MCT-all	15.51	15.50	19.36	21.42	19.69	31.19	32.73	35.36	36.23	34.43
Joint-Straight [22]	10.46	10.37	11.16	14.98	12.86	21.42	22.95	23.37	30.64	24.68
Joint-ASO [23]	10.33	10.31	10.37	15.66	11.67	21.26	21.81	22.08	29.43	23.67
Joint-Grad [24]	9.53	10.23	10.22	15.48	11.55	21.21	21.67	21.91	29.20	23.54
Proposed	9.33	9.77	9.91	14.83	11.86	20.87	21.03	21.23	28.82	22.84
	9.33	9.63	9.81	14.53	11.88	20.78	20.96	21.10	28.41	22.84

TABLE 5. Comparison of speech quality and noise reduction quality scores of 18 models trained using different training approaches on Test-Clean in the noisy LibriSpeech dataset.

Training Approach	Speech Quality					Noise Reduction Quality				
	PESQ	STOI	CN0	CBAK	COVL	SIR	SIR	SAR	SNR	SNR
MCT-noisy	1.7256	0.6867	1.8547	1.1615	1.3937	8.1762	-0.1762	-2.3140	-0.2280	-0.2280
+ standardize-SR	2.8512	0.8271	2.9671	2.5402	2.5430	11.8996	18.4267	12.3713	5.0665	18.5811
Joint-Straight [22]	2.8234	0.8607	2.8330	2.3690	2.2099	8.7618	13.0613	13.1106	4.2454	8.8827
Joint-ASO [23]	2.5478	0.8732	2.8003	2.3241	2.3394	10.6317	16.4708	11.9010	4.5793	10.1932
Joint-Grad [24]	2.5413	0.8739	2.8113	2.3244	2.3466	10.3110	13.1661	12.1499	4.4104	8.9012
Proposed	2.6653	0.8311	3.1284	2.5684	2.4509	10.0398	17.1702	12.4433	4.5182	10.3467

TABLE 6. Comparison of speech quality and noise reduction quality scores of 18 models trained using different training approaches on test-clean in the CHiME-4 test dataset.

Training Approach	Speech Quality					Noise Reduction Quality				
	PESQ	STOI	CN0	CBAK	COVL	SIR	SIR	SAR	SNR	SNR
MCT-noisy	1.8869	0.6106	2.3944	1.7629	1.7216	5.2784	-2.7784	-1.4166	-1.2390	-1.2390
+ standardize-SR	2.8932	0.8991	2.9186	2.4821	2.3711	12.9113	17.9906	12.5413	18.2023	11.5718
Joint-Straight [22]	2.5234	0.8685	2.8130	2.3090	2.2668	11.8877	16.3162	12.1276	9.9932	10.8723
Joint-ASO [23]	2.5459	0.8732	2.6803	2.3241	2.3395	12.3478	16.5440	12.2389	10.9773	11.1873
Joint-Grad [24]	2.5411	0.8789	2.6611	2.3244	2.3496	12.3252	16.6106	12.2722	10.1423	11.1933
Proposed	2.6098	0.8894	3.0731	2.5133	2.4518	12.5444	17.6202	12.5554	10.1864	11.5705

TABLE 8. Ablative study on the effectiveness of different loss combinations in the proposed training approach on ASR performance using Test-Clean in the noisy LibriSpeech dataset (L_{r} is applied to the proposed training approach).

Training Approach	Loss Function			CER (%)					WER (%)				
	L_{r}	L_{e}	L_{d}	Development			Test		Avg.	Development			Avg.
				dev-clean	dev-other	test-clean	test-other			dev-clean	dev-other	test-clean	test-other
MCT-noisy	-	-	-	12.87	12.93	13.04	13.10	12.98	22.77	23.16	22.95	23.47	23.06
	-	-	-	16.83	16.83	16.96	16.90	16.88	28.94	29.03	29.36	29.14	29.12
Proposed	-	-	-	12.70	12.68	12.80	12.92	12.81	22.58	22.76	22.74	23.15	22.81
	-	-	-	11.82	11.28	11.47	11.42	11.38	19.88	20.40	20.18	20.67	20.39

TABLE 9. Ablative study on the effectiveness of various loss combinations in the proposed training approach on speech and noise reduction quality on Test-Clean in the CHiME-4 dataset (L_{r} is applied to the proposed training approach).

Training Approach	Loss Function			Speech Quality					Noise Reduction Quality				
	L_{r}	L_{e}	L_{d}	PESQ	STOI	CN0	CBAK	COVL	SIR	SIR	SAR	SNR	SNR
MCT-noisy	-	-	-	1.7256	0.6867	1.8547	1.1615	1.3937	8.1762	-0.1762	-2.3140	-0.2280	-0.2280
	-	-	-	2.8512	0.8271	2.9671	2.5402	2.5430	11.8996	18.4267	12.3713	5.0665	18.5811
Proposed	-	-	-	2.6653	0.8311	3.1284	2.5684	2.4509	10.0398	17.1702	12.4433	4.5182	10.3467
	-	-	-	2.6653	0.8311	3.1284	2.5684	2.4509	10.0398	17.1702	12.4433	4.5182	10.3467

TABLE 7. Ablative study on the effectiveness of different loss combinations in the proposed training approach on ASR performance on Test-Clean in the CHiME-4 dataset (L_{r} is applied to the proposed training approach).

Training Approach	Loss Function			CER (%)					WER (%)				
	L_{r}	L_{e}	L_{d}	Development			Test		Avg.	Development			Avg.
				dev01-clean	dev01-noisy	dev01-real	test01-clean	test01-real		dev01-clean	dev01-noisy	dev01-real	
MCT-noisy	-	-	-	10.63	10.63	11.27	10.02	12.13	12.13	22.80	23.00	25.57	31.08
	-	-	-	15.51	15.50	19.36	21.42	19.69	19.69	31.19	32.73	35.36	36.23
Proposed	-	-	-	10.33	10.31	10.37	15.66	11.67	11.67	21.26	21.81	22.08	29.43
	-	-	-	9.53	10.23	10.22	15.48	11.55	11.55	21.21	21.67	21.91	29.20
	-	-	-	9.33	9.77	9.91	14.83	11.86	11.86	20.87	21.03	21.23	28.82
	-	-	-	9.33	9.63	9.81	14.53	11.88	11.88	20.78	20.96	21.10	28.41

TABLE 10. Ablative study on the effectiveness of various loss combinations in the proposed training approach on speech and noise reduction quality using test-clean in the CHiME-4 dataset (L_{r} is applied to the proposed training approach).

Training Approach	Loss Function			Speech Quality					Noise Reduction Quality				
	L_{r}	L_{e}	L_{d}	PESQ	STOI	CN0	CBAK	COVL	SIR	SIR	SAR	SNR	SNR
MCT-noisy	-	-	-	1.8869	0.6106	2.3944	1.7629	1.7216	5.2784	-2.7784	-1.4166	-1.2390	-1.2390
	-	-	-	2.8932	0.8991	2.9186	2.4821	2.3711	12.9113	17.9906	12.5413	18.2023	11.5718
Proposed	-	-	-	2.6098	0.8894	3.0731	2.5133	2.4518	12.5444	17.6202	12.5554	10.1864	11.5705
	-	-	-	2.6098	0.8894	3.0731	2.5133	2.4518	12.5444	17.6202	12.5554	10.1864	11.5705

Fig -1: Figure

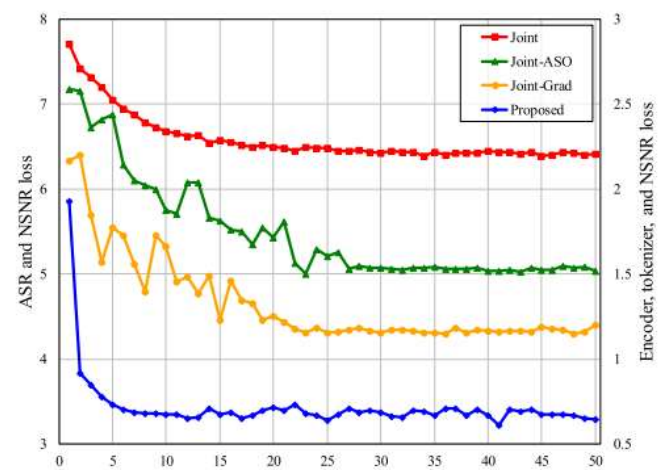


FIGURE 3. Training loss curves on four different training approaches: the left y-axis corresponds to the loss values for Joint, Joint-ASO, Joint-Grad, while the right y-axis to the proposed training approach.

3. CONCLUSIONS

In this study, we proposed a knowledge distillation-based training framework for speech enhancement aimed at improving the robustness of automatic speech recognition (ASR) systems in noisy environments. By leveraging the outputs of a high-capacity

teacher model, a lightweight student model was effectively trained to perform real-time speech enhancement while preserving critical speech features essential for ASR.

The integration of both frame-level and sequence-level loss functions enabled the student model to generalize across a wide range of acoustic conditions, maintaining a balance between noise suppression and speech intelligibility. Experimental evaluations demonstrated that the proposed approach significantly outperforms conventional speech enhancement methods, both in terms of ASR accuracy and computational efficiency.

The resulting student model, with its reduced complexity and latency, is well-suited for deployment in resource-constrained environments such as mobile devices and edge computing platforms. This work highlights the potential of knowledge distillation in advancing noise-robust ASR and sets a foundation for future exploration in multi-task learning, adaptive distillation, and domain generalization in speech technologies.

ACKNOWLEDGEMENT

The author sincerely acknowledges the invaluable guidance, continuous support, and constructive feedback provided by Dr. S. China Venkateswarlu and Dr. V. Siva Nagaraju faculty members of the Department of Electronics and Communication Engineering at the Institute of Aeronautical Engineering (IARE). Their expert advice and encouragement have been instrumental throughout the entire course of this research.

Special thanks are also extended to the faculty and staff of the Institute for providing a conducive academic environment and essential resources that greatly facilitated the successful completion of this work. The author appreciates the support and collaboration of peers and colleagues who contributed their time and expertise.

REFERENCES

- [1] D. Wang and J. Chen, "Supervised speech separation based on deep learning: An overview," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 26, no. 10, pp. 1702–1726, Oct. 2018.
- [2] S. A. Nossier, J. Wall, M. Moniri, C. Glackin, and N. Cannings, "An experimental analysis of deep learning architectures for supervised speech enhancement," *Electronics*, vol. 10, no. 1, p. 17, Dec. 2020.
- [3] J. M. Perero-Codosero, F. M. Espinoza-Cuadros, and L. A. Hernández Gómez, "A comparison of hybrid and end-to-end ASR systems for the IberSpeech-RTVE 2020 speech-to-text transcription challenge," *Appl. Sci.*, vol. 12, no. 2, p. 903, Jan. 2022.
- [4] G. Caldarini, S. Jaf, and K. McGarry, "A literature survey of recent advances in chatbots," *Information*, vol. 13, no. 1, p. 41, Jan. 2022.
- [5] S. E. Eskimez, T. Yoshioka, H. Wang, X. Wang, Z. Chen, and X. Huang, "Personalized speech enhancement: New models and comprehensive evaluation," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, Shenzhen, China, May 2022, pp. 356–360.
- [6] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," in *Proc. Int. Conf. Mach. Learn. (ICML)*, Honolulu, HI, USA, Jul. 2023, pp. 28492–28518.
- [7] G. E. Dahl, D. Yu, L. Deng, and A. Acero, "Context-dependent pre trained deep neural networks for large-vocabulary speech recognition," *IEEE Trans. Audio, Speech, Language Process.*, vol. 20, no. 1, pp. 30–42, Jan. 2012.
- [8] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "Wav2vec 2.0: A framework for self-supervised learning of speech representations," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Dec. 2020, pp. 12449–12460.
- [9] M. L. Seltzer, D. Yu, and Y. Wang, "An investigation of deep neural networks for noise robust speech recognition," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process.*, Vancouver, BC, Canada, May 2013, pp. 7398–7402.
- [10] W. Chan, N. Jaitly, Q. Le, and O. Vinyals, "Listen, attend and spell: A neural network for large vocabulary conversational speech recognition," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, Shanghai, China, Mar. 2016, pp. 4960–4964.
- [11] S. Karita, N. Chen, T. Hayashi, T. Hori, H. Inaguma, Z. Jiang, M. Someki, N. E. Y. Soplin, R. Yamamoto, X. Wang, S. Watanabe, T. Yoshimura, and W. Zhang, "A comparative study on transformer vs RNN in speech applications," in *Proc. IEEE*

Autom. Speech Recognit. Understand. Workshop (ASRU), Dec. 2019, pp. 449–456.

[12] A. Gulati, J. Qin, C. C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu, and R. Pang, “Conformer: Convolution-augmented transformer for speech recognition,” in Proc. Interspeech, Shanghai, China, Oct. 2020, pp. 5036–5040.

[13] W. Han, Z. Zhang, Y. Zhang, J. Yu, C.-C. Chiu, J. Qin, A. Gulati, R. Pang, and Y. Wu, “ContextNet: Improving convolutional neural networks for automatic speech recognition with global context,” in Proc. Interspeech, Shanghai, China, Oct. 2020, pp. 3610–3614.

[14] D. Amodei, S. Ananthanarayanan, R. Anubhai, J. Bai, E. Battenberg, C. Case, J. Casper, B. Catanzaro, Q. Cheng, G. Chen, and J. Chen, “Deep speech 2: End-to-end speech recognition in English and Mandarin,” in Proc. Int. Conf. Mach. Learn. (ICML), New York, NY, USA, Jun. 2016, pp. 173–182.

[15] A. Graves, “Sequence transduction with recurrent neural networks,” 2012, arXiv:1211.3711.

[16] J. Droppo and A. Acero, “Joint discriminative front-end and back-end training for improved speech recognition accuracy,” in Proc. IEEE Int. Conf. Acoust. Speech Signal Process., Toulouse, France, May 2006, pp. 281–284.

[17] K. Kinoshita, T. Ochiai, M. Delcroix, and T. Nakatani, “Improving noise robust automatic speech recognition with single-channel time-domain enhancement network,” in Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP), Barcelona, Spain, May 2020, pp. 7009–7013.

[18] Q. Wang, P. Wang, and D. Wang, “Complex spectral mapping for single and multi-channel speech enhancement and robust ASR,” IEEE/ACM Trans. Audio, Speech, Language Process., vol. 28, pp. 1778–1787, 2020.

[19] Y. Hu, Y. Liu, S. Lv, M. Xing, S. Zhang, Y. Fu, J. Wu, B. Zhang, and L. Xie, “DCCRN: Deep complex convolution recurrent network for phase-aware speech enhancement,” in Proc. Interspeech, Shanghai, China, Oct. 2020, pp. 2472–2476.

[20] K. Shimada, Y. Bando, M. Mimura, K. Itoyama, K. Yoshii, and T. Kawahara, “Unsupervised speech enhancement based on multichannel NMF-informed beamforming for noise-robust automatic speech recognition,” IEEE/ACM Trans. Audio, Speech, Language Process., vol. 27, no. 5, pp. 960–971, May 2019.

[21] B. Schuller, F. Weninger, M. Wöllmer, Y. Sun, and G. Rigoll, “Non negative matrix factorization as noise-robust feature extractor for speech recognition,” in Proc. IEEE Int. Conf. Acoust., Speech Signal Process., Dallas, TX, USA, Mar. 2010, pp. 4562–4565.

[22] L. Li, Y. Kang, Y. Shi, L. Kürzinger, T. Watzel, and G. Rigoll, “Adversarial joint training with self-attention mechanism for robust end-to-end speech recognition,” EURASIP J. Audio, Speech, Music Process., vol. 2021, no. 1, pp. 1–16, Jul. 2021.

[23] D. Ma, N. Hou, V. T. Pham, H. Xu, and E. S. Chng, “Multitask-based joint learning approach to robust ASR for radio communication speech,” in Proc. Asia-Pacific Signal Inf. Process. Assoc. Annu. Summit Conf. (APSIPA ASC), Dec. 2021, pp. 497–502.

[24] Y. Hu, C. Chen, R. Li, Q. Zhu, and E. S. Chng, “Gradient remedy for multi task learning in end-to-end noise-robust speech recognition,” in Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP), Rhodes Island, Greece, Jun. 2023, pp. 1–5.

[25] G. W. Lee and H. K. Kim, “Two-step joint optimization with auxiliary loss function for noise-robust speech recognition,” Sensors, vol. 22, no. 14, p. 5381, Jul. 2022.

[26] A. Pandey, C. Liu, Y. Wang, and Y. Saraf, “Dual application of speech enhancement for automatic speech recognition,” in Proc. IEEE Spoken Lang. Technol. Workshop (SLT), Shenzhen, China, Jan. 2021, pp. 223–228.

[27] C. C. Lee, Y. Tsao, H. M. Wang, and C. S. Chen, “D4AM: A general denoising framework for downstream acoustic models,” in Proc. Int. Conf. Learn. Represent. (ICLR), Kigali, Rwanda, May 2023, pp. 1–17.

[28] M. Yang, J. Konan, D. Bick, Y. Zeng, S. Han, A. Kumar, S. Watanabe, and B. Raj, “Paaploss: A phonetic-aligned acoustic parameter loss for speech enhancement,” in Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP), Rhodes Island, Greece, Jun. 2023, pp. 1–5.

[29] G. Hinton, O. Vinyals, and J. Dean, “Distilling the knowledge in a neural network,” 2015, arXiv:1503.02531.

[30] A. Romero, N. Ballas, S. Ebrahimi Kahou, A. Chassang, C. Gatta, and Y. Bengio, “FitNets: Hints for thin deep nets,” 2014, arXiv:1412.6550.

[31] W. H. Li and H. Bilen, “Knowledge distillation for multi-task learning,” in Proc. Eur. Conf. Comput. Vis. (ECCV), Aug. 2020, pp. 163–176.

[32] G. M. Jacob, V. Agarwal, and B. Stenger, “Online knowledge distillation for multi-task learning,” in Proc. IEEE/CVF Winter Conf. Appl. Comput. Vis. (WACV), Waikoloa, HI, USA, Jan. 2023, pp. 2358–2367.