

InQSS: a speech intelligibility and quality assessment model using a multi-task learning network

Yu-Wen Chen, Yu Tsao

Research Center for Information Technology Innovation, Academia Sinica, Taiwan

Abstract

Speech intelligibility and quality assessment models are essential tools for researchers to evaluate and improve speech processing models. However, only a few studies have investigated multi-task models for intelligibility and quality assessment due to the limitations of available data. In this study, we released TMHINT-QI, the first Chinese speech dataset that records the quality and intelligibility scores of clean, noisy, and enhanced utterances. Then, we propose InQSS, a non-intrusive multi-task learning framework for intelligibility and quality assessment. We evaluated the InQSS on both the training-from-scratch and the pretrained models. The experimental results confirm the effectiveness of the InQSS framework. In addition, the resulting model can predict not only the intelligibility scores but also the quality scores of a speech signal.

Index Terms: intelligibility assessment, quality assessment, self-supervised learning, multi-task neural network

1. Introduction

The speech intelligibility and quality assessment models are important because they provide an efficient way to evaluate the speech processing models, such as the speech enhancement (SE), automatic speech recognition (ASR), and voice conversion models. Various methods have been introduced to measure the speech intelligibility [1, 2, 3] and quality [4, 5, 6]. However, these methods are intrusive measurements that require knowing the corresponding clean speech signals, which are often unavailable in real-world applications. In addition, the results might not correlate well with the listening test results [7, 8, 9].

To address this problem, recent studies have used listening test datasets to build non-intrusive neural network models that can predict human perception. For example, [10, 11, 12, 13] proposed quality assessment models trained on VCC challenge dataset and [9] proposed DNSMOS trained on the DNS challenge dataset [14]. For intelligibility assessment, the models in [15] and [16] were trained using the dataset presented in [17] and a self-collected Japanese corpus, respectively. In addition, [18] used previous quality assessment models as references, and [19] proposed a multi-task model trained on subjective quality scores and objective scores. However, these models can only perform subjective intelligibility or quality assessment but not both. Furthermore, previous datasets are mostly in English, so datasets in other languages are required.

The contributions of this study are as follows:

- We released TMHINT-QI¹, a Chinese speech dataset with subjective quality and intelligibility scores. The dataset includes clean, noisy, and enhanced noisy utterances processed using SE models. To the best of our knowledge, this is the first Chinese dataset that records

the quality and intelligibility scores of enhanced speech signals.

- We propose the InQSS, a multi-task learning framework using training-from-scratch and pretrained self-supervised learning (SSL) models. The experimental results confirmed the effectiveness of incorporating the training of quality and intelligibility predictions in the same model. In addition, the resulting models are the first multi-task intelligibility and quality assessment models trained on the subjective quality and intelligibility scores.
- We tested and demonstrated the effectiveness of using scattering coefficients [20], which are helpful for several signal processing tasks [21, 22, 23] but have not yet been tested in a speech assessment model.

2. TMHINT-QI dataset

The TMHINT-QI dataset is a Chinese speech dataset with subjective quality and intelligibility scores. TMHINT-QI includes clean utterances recorded in a quiet environment, artificially contaminated noisy utterances, and enhanced noisy utterances processed by the SE models. The goal of releasing TMHINT-QI is to facilitate research on speech quality and intelligibility assessment models, which can later be used to assess the performance of speech processing, thereby guiding researchers to develop models that can generate results with better human perception.

2.1. Data preparation

The TMHINT-QI dataset used the TMHINT sentences [24] and was recorded in a 16-bit format at a 16-kHz sampling rate. Each utterance contained 10 Chinese characters of approximately 3 s in duration and was recorded in a quiet room. The recorded utterances were divided into two parts.

The first part consisted of 3 female and 3 male speakers, each reading 200 sentences with a total of 1200 clean utterances. To form noisy utterances, each clean utterance was contaminated with five randomly sampled noise types from a 100-noise type dataset [25] at 8 different SNR levels (± 1 dB, ± 4 dB, ± 7 dB, and ± 10 dB). These clean-noisy paired utterances were then further used as training data for the three neural-network-based SE models including FCN [26], DDAE [27], and transformer-based SE [28] (denoted as Trans).

The second part of the recorded TMHINT utterances contained one female and one male speaker. Each speaker recorded 115 sentences, a total of 230 clean utterances were recorded. The second part of the utterances was used for the formal listening test. The recorded clean utterances were contaminated with four types of noise (babble, street, pink, and white) at four SNR levels (-2, 0, 2, and 5) to form the noisy speech. These

¹Download TMHINT-QI: <https://github.com/yuwchen/InQSS>

noisy speeches were then processed using five SE models (KLT [29], MMSE [30], FCN, DDAE, and Trans). Finally, the clean, noisy, and enhanced utterances were combined to form the listening test utterance pool.

2.2. Listening test

In the listening test, the participants scored the utterances based on quality and intelligibility. The quality score is on a scale of 1-5, where 1 represents the lowest perceived quality and 5 indicates the highest perceived quality. The intelligibility score is the number of characters a participant can recognize from an utterance. Because each TMHINT sentence contains 10 characters, the intelligibility score is within the range of 0-10, where the score represents the number of characters a participant can recognize in the sentence.

The listening test was divided into pretest and formal test. In the pretest, the participants had to listen to and score five clean utterances from the SE models' training data. We informed the participants that these clean utterances had a quality score of 5, so the upper bound of quality scores might be less affected by the headphones they used. In addition, participants' ability to recognize all characters of clean utterances provided an indicator for evaluating their understanding of Chinese and hearing status.

During the formal test, the participants listened to and scored 103 utterances chosen from the listening test utterance pool. The samples in the formal test started with two randomly chosen files, and the other 101 files included 5 clean utterances and 96 noisy or enhanced utterances. To avoid the problem of a data imbalance in later studies, the number of testing files was designed to be 96 such that each participant would score the same number of files under different SNRs, methods, or noise types. For example, because there are six processing methods (five SE methods and one without processing), a participant listened to 16 (96/6) samples for each method.

In total, TMHINT-QI contains 24,408 samples with 14,919 unique utterances (including example files). These data were collected from 226 participants (94 males and 132 females) between 20 and 50 years of age.

2.3. Characteristics of TMHINT-QI

2.3.1. Comparison between objective assessment metrics and the listening test

We compared two objective assessment metrics (PESQ [4] for quality and STOI [1] for intelligibility) with the listening test results. Figure 1 (left) shows the average PESQ, STOI, quality scores, and intelligibility scores of the SE methods under different SNRs. For PESQ and STOI, the results in (a) and (c) show that noisy utterances without SE (w/o SE) achieved the worst performance, and neural network-based SE methods (FCN, DDAE, and Trans) outperformed the traditional SE methods (KLT and MMSE). These results are also in line with previous studies and our expectations. However, unlike the PESQ and STOI results, the listening results in (b) and (d) show that noisy utterances without applying SE achieved the highest scores in both quality and intelligibility.

The results of our listening test align with the finding of previous studies [7, 9, 8] that the objective assessment metrics might not correlate well with human perception. Note that in the listening test results, the higher the SNRs were, the higher the quality and intelligibility scores the utterances received. This result matches the intuition; thus, we believe the collected

scores can reflect general perceptions to some extent.

2.3.2. Characteristics of quality and intelligibility scores

Figure 1 (right) shows the histograms of the quality and the intelligibility scores, which exclude the scores of clean utterances. The quality scores were distributed close to a Gaussian with a mean of approximately 3. The distribution of intelligibility scores had a strong left skewness and a peak at the maximum. This result indicates that the participants can recognize whole sentences in most samples. In addition, we found similar distributions of the quality and intelligibility scores in the VCC 2018 [10] and ADFD datasets [31], respectively.

2.3.3. Correlation between quality and intelligibility scores

Figure 2 shows the correlation between intelligibility and quality scores of testing data in 4.1. The high correlation between the intelligibility and quality scores motivates the study of using a multi-task learning network to enhance the performance of a single task.

3. Proposed InQSS

3.1. Related works

3.1.1. Scattering transform

The scattering transform was proposed in [20] for building robust, time-shift-invariant, and informative signal features. Previous studies have shown that such features can be successfully applied in several signal processing tasks [21, 22, 23]. However, scattering coefficients have not yet been tested in speech assessment.

3.1.2. MOSNet

MOSNet [10] is a non-intrusive quality assessment model based on convolutional neural network (CNN) and bidirectional long short-term memory (BLSTM). In addition, MOSNet uses a combination of utterance- and frame-level losses such that the prediction is more correlated with human ratings [32].

3.1.3. SSL

SSL aims to generate a general representation of the input signal [33]. Previous studies have demonstrated the strong performance of using SSL features as additional inputs [13] or fine-tuning SSL models for various tasks [34, 35].

3.2. InQSS framework

We propose InQSS, a multi-task learning framework for Intelligibility and Quality aSSessment. We tested the framework on three model structures, including a model trained from scratch (InQSS-MOSNet), a model fine-tuned from a pretrained SSL (InQSS-SSL), and an ensemble model (InQSS-MOSSSL).

The InQSS-MOSNet is based on the MOSNet. We use the same utterance- and frame-wise losses and a similar CNN-BLSTM structure. We improve the MOSNet by adding a multi-task structure and incorporating the scattering coefficients as additional input features. The input spectrogram and scattering coefficients, which are a concatenation of first- and second-order scattering coefficients, first pass through two separated CNNs and are then concatenated as the input of the BLSTMs. Finally, the outputs of the BLSTMs go through the dense layers and give the predicted scores. The InQSS-MOSNet is trained

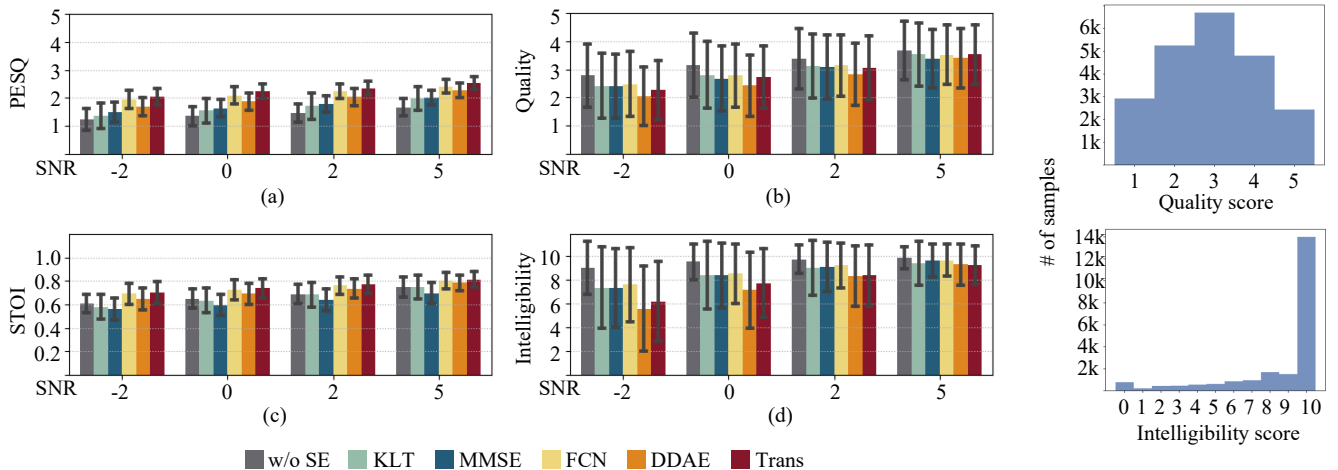


Figure 1: Comparison between objective assessment metrics and the listening test (left), and histograms of the quality scores and the intelligibility scores (right).

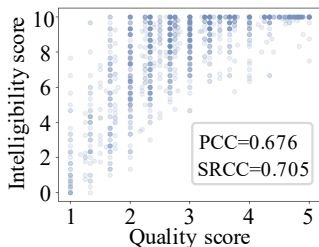


Figure 2: Scatter plot of intelligibility scores versus quality scores. Each dot represents a sample point. The darker the color, the more samples that are gathered.

with the L2 loss of intelligibility and quality scores. The model structure of InQSS-MOSNet is shown in Figure 3.

For the InQSS-SSL, we fine-tune a pretrained SSL model by average-pooling the model’s output embeddings and adding a dense output layer for intelligibility and quality prediction. Then, the InQSS-SSL is trained with the L1 loss of intelligibility and quality scores. The model structure of InQSS-SSL is shown in Figure 4. In InQSS-MOSNet and InQSS-SSL, both the prediction paths are trained simultaneously. For the InQSS-MOSSSL, we first finish the training of the InQSS-MOSNet and InQSS-SSL separately. Then, we average the predictions as final results.

4. Experiments

4.1. Experimental settings

In this study, we used the TMHINI-QI dataset to evaluate the proposed assessment models. To reduce the effect of the bias of the participants on the assessment, we collected utterances with at least three sampled scores as testing data. Then, we used the average scores as target scores. In total, the testing set contained 1978 unique utterances. The remaining 17,448 samples in TMHINI-QI were used for training. Note that the example files in TMHINI-QI were excluded from the testing data and only used for training. During training, we randomly sampled 10% of the data for validation and used these validation data as a reference for early stopping and model selection. The model

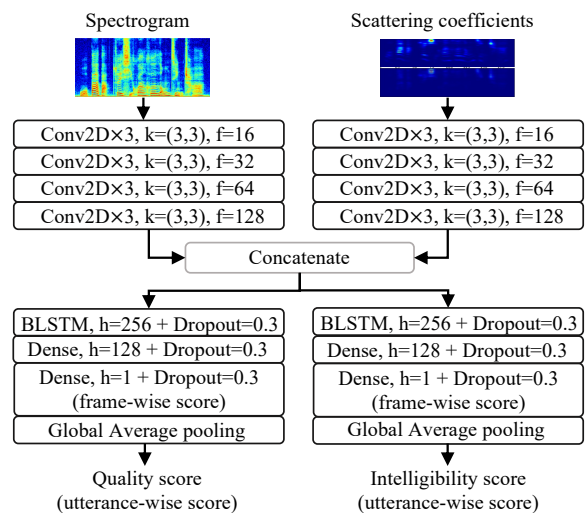


Figure 3: Model structure of InQSS-MOSNet. In convolution layers, k and f denote the kernel size and the number of filters, respectively. In the BLSTMs and dense layers, h denotes the hidden size. Note that a dense layer is a time-distributed dense layer, and the ReLU activation function following every convolution and dense layer is not shown in the figure.

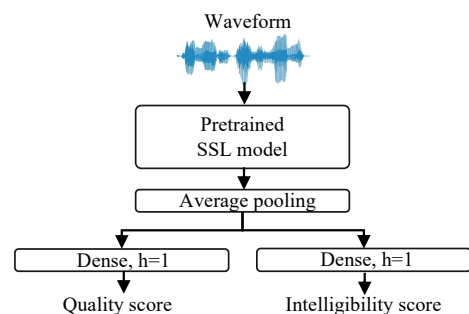


Figure 4: Model structure of InQSS-SSL. In this study, the SSL model we used is w2v_small in the Fairseq toolkit [36].

Table 1: Performance of different intelligibility assessment methods. Here, scat, spec, wav are abbreviations of the scattering coefficients, spectrogram, and raw waveform, respectively.

Model / Method	Input	MSE	PCC	SRCC
In-MOSNet-S _I	spec	2.562	0.695	0.610
In-MOSNet-S _{II}	scat	2.425	0.708	0.633
In-MOSNet-S _{III}	spec+scat	2.393	0.714	0.642
InQSS-MOSNet	spec+scat	2.117	0.755	0.682
In-SSL	wav	2.571	0.749	0.645
InQSS-SSL	wav	2.552	0.754	0.664
In-MOSSSL	wav	2.015	0.777	0.668
InQSS-MOSSSL	spec+scat	2.017	0.791	0.700
	wav			
STOI [1]	-	5.573	0.482	0.461
Google-ASR	-	7.305	0.710	0.679

performances in the following are the average performance of four models trained with different training and validation splits.

The spectrogram was calculated using the short-time Fourier transform with a window length of 512 and a hop length of 256 samples. We used the Kymatio [37] toolkit to conduct scattering transform. The scattering coefficients were calculated by setting the resolution of the first-order wavelet index set Q_1 , and averaging scale of the low-pass averaging filter J , to 8 each.

In addition, we applied sample-wise min-max normalization to the input spectrogram and scattering coefficients and rescaled the intelligibility scores to 0-5 during training such that both the input features and output targets have similar scales. Finally, we evaluated our results using the mean square error (MSE), Pearson’s correlation coefficient (PCC), and Spearman’s rank correlation coefficient (SRCC).

4.2. Experimental results

4.2.1. Evaluation of intelligibility assessment

First, we conducted an ablation study to test the performance of incorporating scattering coefficients. The results of In-MOSNet-S_{II} and In-MOSNet-S_I in Table 1 show that scattering coefficients are more useful for intelligibility prediction than spectrograms. Also, the combination of scattering coefficients and spectrograms performs better than using only spectrograms or only scattering coefficients.

Then, we tested the performance of incorporating quality scores into the intelligibility assessment. In Table 1, an *In*-system has a similar model structure as an *InQSS*-system but does not have a path to predict the quality scores. Comparing the performance of *InQSS*-systems with *In*-systems, the results show that the information regarding the quality scores can improve the performance of an intelligibility assessment. In addition, the ensemble model InQSS-MOSSSL achieves the best performance within the tested models.

Finally, we show the performance of STOI [1] and Google-ASR [38] on the TMHINT-QI dataset. The Google-ASR results were calculated by computing the distance between the predicted and ground-truth sentences. The results in Table 1 indicate the ASR results have a high correlation with the listening test results, whereas the STOI scores are less consistent with the listening test.

Table 2: Performance of different quality assessment models. In the model column, the star mark * indicates whether the model has access to the training set of the dataset.

Model	Dataset	MSE	PCC	SRCC
Q-MOSNet*	TMHINT-QI	0.439	0.753	0.698
InQSS-MOSNet*	TMHINT-QI	0.422	0.763	0.715
Q-SSL*	TMHINT-QI	0.388	0.794	0.750
InQSS-SSL*	TMHINT-QI	0.365	0.800	0.754
InQSS-MOSSSL*	TMHINT-QI	0.353	0.804	0.759
DNSMOS [9]	TMHINT-QI	0.915	0.496	0.311
NISQA [40]	TMHINT-QI	3.140	0.529	0.348
SSL [35]	TMHINT-QI	4.417	0.574	0.405
Q-MOSSSL	DAPS	1.261	0.617	0.599
InQSS-MOSSSL	DAPS	1.100	0.639	0.639
DNSMOS [9]	DAPS	0.665	0.515	0.510
NISQA [40]	DAPS	0.663	0.519	0.389
SSL [35]	DAPS	0.475	0.710	0.718

4.2.2. Evaluation of quality assessment

We tested the performance of incorporating intelligibility scores into the quality assessment. In Table 2, a *Q*-system has a similar model structure as an *InQSS*-system but does not have a path to predict the intelligibility scores. Compared the performance of *InQSS*-systems with *Q*-systems, the results indicate that the information regarding the intelligibility scores can improve the performance of the quality assessment. In addition, we notice that models trained on other datasets do not generalize well to the TMHINT-QI dataset.

We then tested the model’s performance on the DAPS dataset [39]. The results listed in Table 2 show that InQSS-MOSSSL performs better than Q-MOSSSL. However, the SSL model in [35] obtains the best performance. The reason might be that the SSL model in [35] was trained on a much larger speech quality dataset than the TMHINT-QI, and therefore has a better generalizability than our model.

In addition, we find that MSE evaluation results are inconsistent with the PCC and SRCC results on out-of-domain datasets. One possible reason is that different datasets might have similar clean utterances but have diverse noisy utterances processed by different methods. Therefore, because the quality scores of noisy utterances in different datasets have discrepant scoring standards, models are less capable of predicting exact scores for out-of-domain data.

5. Conclusions

In this study, we collected and analyzed the subjective and objective intelligibility and quality scores of clean, noisy, and enhanced utterances. Then, we released the dataset named TMHINT-QI. Moreover, we propose InQSS, a non-intrusive multi-task learning framework for intelligibility and quality assessment. The experimental results demonstrate that a multi-task learning network can improve the performance of a single task without increasing the model complexity. In addition, SSL-based models can achieve high performance on multi-task speech assessment and require less time to convergence than the training-from-scratch models. Finally, a simple ensemble approach, averaging the final predictions of two models, can effectively improve the results.

6. References

- [1] C. H. Taal, R. C. Hendriks, R. Heusdens, and J. Jensen, "An algorithm for intelligibility prediction of time-frequency weighted noisy speech," *IEEE/ACM Trans. Audio Speech Lang.*, vol. 19, no. 7, pp. 2125–2136, 2011.
- [2] S. Jørgensen, S. D. Ewert, and T. Dau, "A multi-resolution envelope-power based model for speech intelligibility," *J. Acoust. Soc. Am.*, vol. 134, no. 1, pp. 436–446, 2013.
- [3] S. Van Kuyk, W. B. Kleijn, and R. C. Hendriks, "An instrumental intelligibility metric based on information theory," *IEEE Signal Process. Lett.*, vol. 25, no. 1, pp. 115–119, 2017.
- [4] A. W. Rix, J. G. Beerends, M. P. Hollier, and A. P. Hekstra, "Perceptual evaluation of speech quality (PESQ)—a new method for speech quality assessment of telephone networks and codecs," in *Proc. ICASSP 2001*.
- [5] R. Huber and B. Kollmeier, "PEMO-Q—A new method for objective audio quality assessment using a model of auditory perception," *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 14, no. 6, pp. 1902–1911, 2006.
- [6] A. Hines, J. Skoglund, A. Kokaram, and N. Harte, "ViSQOL: an objective speech quality model," *EURASIP*, vol. 2015 (13), pp. 1–18, 2015.
- [7] C. K. Reddy, E. Beyrami, J. Pool, R. Cutler, S. Srinivasan, and J. Gehrke, "A scalable noisy speech dataset and online subjective test framework," in *Proc. ICASSP 2019*.
- [8] H. Li and J. Yamagishi, "Noise tokens: Learning neural noise templates for environment-aware speech enhancement," *arXiv preprint arXiv:2004.04001*, 2020.
- [9] C. K. Reddy, V. Gopal, and R. Cutler, "DNSMOS: A non-intrusive perceptual objective speech quality metric to evaluate noise suppressors," in *Proc. ICASSP 2021*.
- [10] C.-C. Lo, S.-W. Fu, W.-C. Huang, X. Wang, J. Yamagishi, Y. Tsao, and H.-M. Wang, "MOSNet: deep learning based objective assessment for voice conversion," in *Proc. INTERSPEECH 2019*.
- [11] Y. Leng, X. Tan, S. Zhao, F. Soong, X.-Y. Li, and T. Qin, "MBNet: Mos prediction for synthesized speech with mean-bias network," in *Proc. ICASSP 2021*.
- [12] Y. Choi, Y. Jung, and H. Kim, "Neural MOS prediction for synthesized speech using multi-task learning with spoofing detection and spoofing type classification," in *Proc. SLT 2021*.
- [13] W.-C. Tseng, C.-y. Huang, W.-T. Kao, Y. Y. Lin, and H.-y. Lee, "Utilizing self-supervised representations for MOS prediction," in *Proc. INTERSPEECH 2021*.
- [14] C. K. Reddy, V. Gopal, R. Cutler, E. Beyrami, R. Cheng, H. Dubey, S. Matussevych, R. Aichner, A. Aazami, S. Braun *et al.*, "The interspeech 2020 deep noise suppression challenge: Datasets, subjective testing framework, and challenge results," *arXiv preprint arXiv:2005.13981*, 2020.
- [15] C. Spille, S. D. Ewert, B. Kollmeier, and B. T. Meyer, "Predicting speech intelligibility with deep neural networks," *Computer Speech & Language*, vol. 48, pp. 51–66, 2018.
- [16] K. Arai, S. Araki, A. Ogawa, K. Kinoshita, T. Nakatani, K. Yamamoto, and T. Irino, "Predicting speech intelligibility of enhanced speech using phone accuracy of DNN-based ASR system," in *Proc. INTERSPEECH 2019*.
- [17] W. Schubotz, T. Brand, B. Kollmeier, and S. D. Ewert, "Monaural speech intelligibility and detection in maskers with varying amounts of spectro-temporal speech features," *J. Acoust. Soc. Am.*, vol. 140, no. 1, pp. 524–540, 2016.
- [18] P. Manocha, B. Xu, and A. Kumar, "NORESQA: a framework for speech quality assessment using non-matching references," in *Proc. NeurIPS 2021*.
- [19] Z. Zhang, P. Vyas, X. Dong, and D. S. Williamson, "An end-to-end non-intrusive model for subjective and objective real-world speech assessment using a multi-task framework," in *Proc. ICASSP 2021*.
- [20] S. Mallat, "Group invariant scattering," *Commun. Pure Appl. Math.*, vol. 65, no. 10, pp. 1331–1398, 2012.
- [21] J. Andén and S. Mallat, "Deep scattering spectrum," *IEEE Trans. Signal Process.*, vol. 62, no. 16, pp. 4114–4128, 2014.
- [22] N. Zeghidour, G. Synnaeve, M. Versteegh, and E. Dupoux, "A deep scattering spectrum—deep Siamese network pipeline for unsupervised acoustic modeling," in *Proc. ICASSP 2016*.
- [23] W. Ghezaiel, L. Brun, and O. Lézoray, "Hybrid network for end-to-end text-independent speaker identification," in *Proc. ICPR 2021*.
- [24] M. Huang, "Development of Taiwan Mandarin hearing in noise test," *Department of speech language pathology and audiology, National Taipei University of Nursing and Health science*, 2005.
- [25] G. Hu and D. Wang, "A tandem algorithm for pitch estimation and voiced speech segregation," *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 18, no. 8, pp. 2067–2079, 2010.
- [26] S.-W. Fu, Y. Tsao, X. Lu, and H. Kawai, "Raw waveform-based speech enhancement by fully convolutional networks," in *Proc. APSIPA 2017*.
- [27] X. Lu, Y. Tsao, S. Matsuda, and C. Hori, "Speech enhancement based on deep denoising autoencoder," in *Proc. INTERSPEECH 2013*.
- [28] J. Kim, M. El-Khamy, and J. Lee, "T-gsa: Transformer with gaussian-weighted self-attention for speech enhancement," in *Proc. ICASSP 2020*.
- [29] A. Rezayee and S. Gazor, "An adaptive KLT approach for speech enhancement," *IEEE Trans. Speech and Audio Process.*, vol. 9, no. 2, pp. 87–95, 2001.
- [30] Y. Ephraim and D. Malah, "Speech enhancement using a minimum mean-square error log-spectral amplitude estimator," *IEEE Trans Acoust.*, vol. 33, no. 2, pp. 443–445, 1985.
- [31] A. H. Andersen, J. M. De Haan, Z.-H. Tan, and J. Jensen, "Non-intrusive speech intelligibility prediction using convolutional neural networks," *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 26, no. 10, pp. 1925–1939, 2018.
- [32] S.-W. Fu, Y. Tsao, H.-T. Hwang, and H.-M. Wang, "Quality-Net: An end-to-end non-intrusive speech quality assessment model based on BLSTM," in *Proc. INTERSPEECH 2018*.
- [33] S. Liu, A. Mallol-Ragolta, E. Parada-Cabeleiro, K. Qian, X. Jing, A. Kathan, B. Hu, and B. W. Schuller, "Audio self-supervised learning: A survey," *arXiv preprint arXiv:2203.01205*, 2022.
- [34] S. Siriwardhana, A. Reis, R. Weerasekera, and S. Nanayakkara, "Jointly fine-tuning "BERT-like" self supervised models to improve multimodal speech emotion recognition," in *Proc. INTERSPEECH 2020*.
- [35] E. Cooper, W.-C. Huang, T. Toda, and J. Yamagishi, "Generalization ability of MOS prediction networks," in *Proc. ICASSP 2022*.
- [36] M. Ott, S. Edunov, A. Baevski, A. Fan, S. Gross, N. Ng, D. Grangier, and M. Auli, "fairseq: a fast, extensible toolkit for sequence modeling," in *Proc. NAACL-HLT 2019*.
- [37] M. Andreux, T. Angles, G. Exarchakis, R. Leonarduzzi, G. Rochette, L. Thiry, J. Zarka, S. Mallat, J. Andén, E. Belilovsky *et al.*, "Kymatio: scattering transforms in Python," *J. Mach. Learn. Res.*, vol. 21, no. 60, pp. 1–6, 2020.
- [38] A. Zhang, "Speech recognition (version 3.8)," 2017 (accessed on October 06, 2021), https://github.com/Uberi/speech_recognition#readme.
- [39] J. Su, A. Finkelstein, and Z. Jin, "Perceptually-motivated environment-specific speech enhancement," in *Proc. ICASSP 2019*.
- [40] G. Mittag, B. Naderi, A. Chehadi, and S. Möller, "NISQA: A deep CNN-self-attention model for multidimensional speech quality prediction with crowdsourced datasets," in *Proc. INTERSPEECH 2021*.