# Improving Speech Enhancement through Fine-Grained Speech Characteristics

*Muqiao Yang[1†], Joseph Konan[1†], David Bick[1†], Anurag Kumar[2], Shinji Watanabe[1], Bhiksha Raj[1]*

[1] Carnegie Mellon University
[2] Meta Reality Labs Research

{muqiaoy, jkonan, dbick, bhiksha}@cs.cmu.edu, swatanab@andrew.cmu.edu, anuragkr@fb.com

## Abstract

While deep learning based speech enhancement systems have made rapid progress in improving the quality of speech signals, they can still produce outputs that contain artifacts and can sound unnatural. We propose a novel approach to speech enhancement aimed at improving perceptual quality and naturalness of enhanced signals by optimizing for key characteristics of speech. We first identify key acoustic parameters that have been found to correlate well with voice quality (e.g. jitter, shimmer, and spectral flux) and then propose objective functions which are aimed at reducing the difference between clean speech and enhanced speech with respect to these features. The full set of acoustic features is the extended Geneva Acoustic Parameter Set (eGeMAPS), which includes 25 different attributes associated with perception of speech. Given the non-differentiable nature of these feature computation, we first build differentiable estimators of the eGeMAPS and then use them to fine-tune existing speech enhancement systems. Our approach is generic and can be applied to any existing deep learning based enhancement systems to further improve the enhanced speech signals. Experimental results conducted on the Deep Noise Suppression (DNS) Challenge dataset shows that our approach can improve the state-of-the-art deep learning based enhancement systems.

**Index Terms**: speech enhancement, eGeMAPS, acoustic parameters, perceptual quality, explainable enhancement evaluation.

## 1. Introduction

Speech enhancement (SE) is aimed at improving the quality and intelligibility of speech signals that have been degraded by noise signals, reverberations and other factors. The need for speech enhancement arises in a variety of applications; communications through VoIP, mobile phones, in hearing aids, in downstream applications such as automatic speech recognition, and many others. A large body of SE literature focuses on single-channel speech enhancement and our focus in this paper is also on single-channel SE.

Classical approaches for SE have been signal processing driven. These include methods such as spectral subtraction and Wiener filtering [1, 2]. However, much of the recent progress in SE comes from deep learning methods. A variety of deep neural network (DNN) based speech enhancement methods have been proposed and improvements have been made through every subsequent architecture [3, 4, 5, 6, 7, 8]. A key component of these approaches is the loss function used for training the neural networks. Once again a wide variety of methods have been proposed, but generally these loss functions try to minimize $L_1/L_2$ loss between the network output and the target clean speech, sometimes in time-domain [9], sometimes in

time-frequency domain [10]. However, these point-wise differences fail to capture key speech characteristics which can lead to enhanced speech with artifacts and unnatural speech [11]. For example, these flaws can arise due to failure to capture pitch differences [12] and ignoring low-energy regions that represent phonemes such as fricatives or plosives [13].

The aforementioned problems led to recent speech enhancement research that specifically optimizes for perceptual quality of the signals. One branch of these works uses differentiable estimators of existing perceptual metrics to optimize their model [14, 15, 16], including Perceptual Evaluation of Speech Quality (PESQ) [17] and Short-Time Objective Intelligibility (STOI) [18]. However, these objective metrics correlate well with human perception only to a certain degree [19, 11, 20] and the benefits of optimizing these metrics are often observed to be limited [21, 15, 22].

Another class of works in this direction does not target any specific metric but attempts to capture perceptual quality through auxiliary losses. Common approaches involve using differences between filter banks or deep representations [23, 24, 25]. However, these losses also provide only implicit supervision and do not specifically target speech characteristics, leading to limited improvements.

In this paper, we argue that optimizing many specific acoustic qualities - the fine-grained features - of enhanced speech will result in improved perceptual quality. These features are important because they have shown to correlate with perceptual quality, such as in breathy or harsh voice quality [26, 27], and raspy voice quality [28]. They have even shown correlation with perception of a speaker's personality, such as warmth [29]. We show that two state-of-the-art models produced enhanced speech with different values for these features than clean speech. This shows the potential of our method compared to metric or other loss optimization.

Unlike any prior works, we explicitly optimize for these important acoustic properties and force the network to pay attention to these low-level features. We use the extended Geneva Minimal Parameter Set (eGeMAPS) [30] as the full feature set, for its coverage of many categories of speech information. Every feature was selected for its association with cues of voice quality or emotion. Existing calculations for these features are non-differentiable, so we train an eGeMAPS prediction network that can plug into other learning pipelines. Its effectiveness is demonstrated by mean-squared error (MSE) of predicted and ground-truth eGeMAPS. After 50 epochs, we reach MSE of $1.9 \cdot 10^{-4}$ on train and $2.6 \cdot 10^{-4}$ on test.

We verify our hypothesis experimentally, by fine-tuning two existing speech enhancement systems with our auxiliary loss to minimize the difference between eGeMAPS features for enhanced and clean speech. Specifically, we chose two mod-

---

† Equal contribution. Reverse alphabetical order.

els that rate highly on standard metrics, FullSubNet[1] [31] and Demucs[2] [9]. Both were trained on data from the Deep Noise Suppression (DNS) Challenge of InterSpeech 2020 [32], and provided pre-trained models. We initialize from these checkpoints, and fine-tune on this dataset because of its large size and variety of noise-types. Our experimental results on these systems show improvements in widely used metrics, PESQ, STOI and MOS (DNSMOS [19] as proxy for MOS).

# 2. Related Work

Our approach falls into the space of speech enhancement methods that are aimed for improving perceptual quality. There are two primary classes of works in this area: (a) objective metric optimization, and (b) feature loss optimization.

## 2.1. Objective Metric Optimization

One recent approach targeting a metric used the generative adversarial network (GAN) framework. The discriminator is a differentiable PESQ estimator, and the generator is the enhancement model [33, 14]. Other approaches use the metrics as rewards within a reinforcement learning context [15], or formulate a differentiable equation that approximates the metric [16]. However, it is hard to obtain strong improvements through implicit supervision from a collapsed metric. Similar to these other methods, we use a neural network as a differentiable estimator of previously non-differentiable quantities. However, unlike these works, we explicitly optimize for natural speech characteristics captured through a set of low-level acoustic features.

## 2.2. Feature Losses

This class of works cover a broader spectrum but they generally try to incorporate domain-knowledge in forms of feature-based loss functions. For example, [23] attempts to bring domain-knowledge from computational models of human auditory perception by minimizing the difference between filter-banks meant to stimulate human cochlea [23]. [13] encodes linguistic domain-knowledge by matching senone logits between clean and enhanced speech. [25] injects emotion and speaker related information. Some have also argued for using deep speech representations for improving enhancement models [24]. Among these methods, our method is along the lines of the first two approaches, as we also target specific features of speech. However, our approach is very different from these as we rely on differentiable estimators of well-defined acoustic features in eGeMAPS set to capture speech characteristics.

# 3. Proposed Framework

## 3.1. Acoustic Parameters

The eGeMAPS is a collection of acoustic parameters that provide details about an audio signal [30]. For example, consider *jitter* features. A perfect sinusoid has a consistent period, but the period length for a speech signal varies by milliseconds (ms) from period to period. The difference between these durations for the fundamental frequency (F0) is defined as *jitter*. Similarly, the amplitude reaches slightly different heights for each period. Shimmer is the difference between peak amplitudes for consecutive periods of F0.

---

[1] https://github.com/haoxiangsnr/FullSubNet
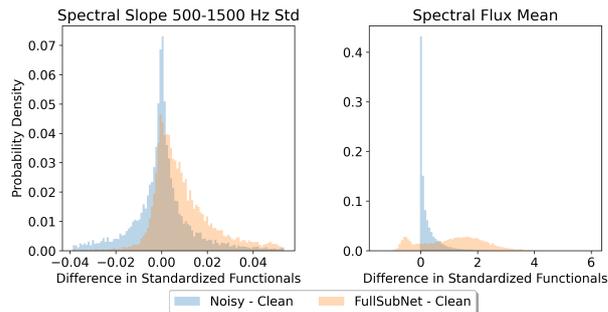[2] https://github.com/facebookresearch/denoiser



Figure 1: *Example distributions of differences between standardized functionals. In both plots we see differences between enhanced and clean speech, with the right plot showing greater differences than noisy speech.*
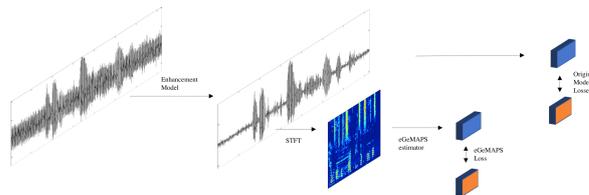


Figure 2: *Our fine-tuning pipeline for arbitrary model. Our proposed loss is on the left. Note that the original model losses can be one or many, and in time domain, frequency domain, or other.*

Researchers have discovered that jitter and shimmer correlate with breathy or harsh voice quality [26, 28]. Another acoustic parameter, harmonics-to-noise ratio, correlates well with raspy vocal quality [27]. Spectral flux and spectral slope have been correlated with the perception of an individual's personality, such as their warmth [29]. Each of the 25 eGeMAPS were chosen by speech experts based on literature associating them with the perception of voice [30].

Given this association of acoustic parameters and perception of voice, we hypothesize that the eGeMAPS could help us quantify the perceptual quality of enhanced speech. We also hypothesize that inducing similarities in these features could improve the quality of the enhancement.

The eGeMAPS are produced in a two-step process. Initially, the 25 features are calculated for windows of 30 ms, with a stride of 10 ms. These are called low-level descriptors (LLDs). Next, statistics are calculated over LLDs, including mean, standard deviation, and percentiles. These are called functionals, and eGeMAPS defines 88 functionals per input. Let $D = 88$ denote the size of each eGeMAPS functionals vector for an utterance.

## 3.2. Visualization of eGeMAPS Differences

First, we show that enhanced speech differs from clean across many features. In Figure 1 we plot the differences in eGeMAPS functionals of clean speech with both enhanced and noisy speech [30]. Since the functionals have different scales, to accurately reflect the scale of differences in the plots, we standardize each functional according to the mean and standard deviation (std) of the *clean* functional values.

For each utterance in the data, indexed by $n \in 1, \ldots, N$, we have a vector of differences $\mathbf{d}^{(n)} \in \mathbb{R}^D$. Each $\mathbf{d}_i^{(n)}, i \in \{1, \ldots, D\}$ is the difference between clean and enhanced/noisy
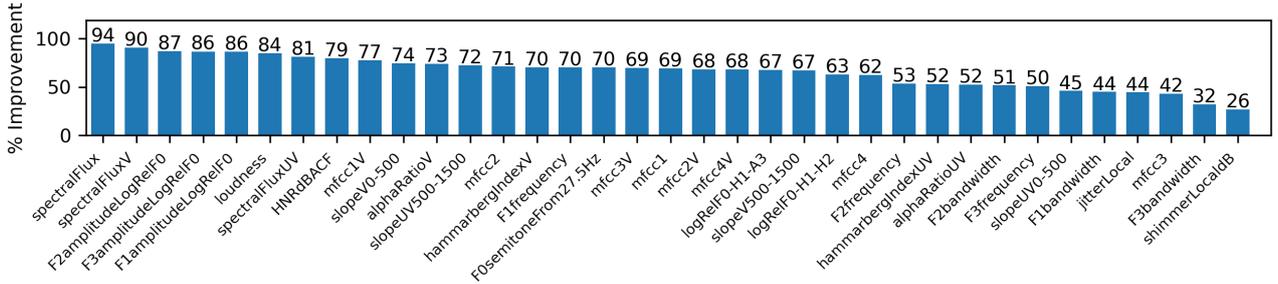
Figure 3: *Percent improvement of the eGeMAPS functionals before and after our fine-tuning. Percent improvement is measured in reduction of difference from clean speech. Labels of openSMILE[34] features are included to allow searching of feature names.*

for eGeMAPS functional $i$. Figure 1 shows $\mathbf{d}^{(n)}$ $\forall n \in \{1, \dots, N\}$ for two functionals: the standard deviation of spectral slope and the mean of spectral flux. The x-axis is the difference in standardized values. The y-axis is the probability of an observation falling in a bin of x-values.

The ideal distribution for eGeMAPS differences of enhanced and clean speech would look close to a Dirac delta function with all mass centered at 0. It would show enhanced speech with matching acoustic features to clean, which we hypothesize provides high perceptual quality. When evaluating the model, any mean different from 0 shows worse performance, and higher variance is worse because it shows more differences of higher values.

In Figure 1, we see that enhanced speech differs from clean speech in a substantial way. The right-hand side shows that enhanced speech differs even more than noisy speech. These large differences show potential for our optimization to induce closer feature values for enhanced and clean speech and improve perceptual quality and naturalness of speech.

### 3.3. Modeling

We now describe our modeling approach to reduce the difference between eGeMAPS for enhanced and clean speech. We first describe each model's inputs and outputs, and how these plug into our method.

Demucs takes in the noisy time-domain signal $\mathbf{x} \in \mathbb{R}^L, \mathbf{x} = \mathbf{y} + \mathbf{n}$, where $\mathbf{y} \in \mathbb{R}^L$ is the clean waveform and $\mathbf{n} \in \mathbb{R}^L$ is the noise waveform. It directly outputs the enhanced waveform $\hat{\mathbf{y}} \in \mathbb{R}^L, \hat{\mathbf{y}} = g_\theta(\mathbf{x})$. Let $g$ be the speech enhancement model parameterized by weights $\theta$, in this case Demucs. The weights are trained to minimize the $L_1$ loss between clean and enhanced waveforms. It also supplements with a loss between spectrograms, calculated by the short-time Fourier transform (STFT). For further details on the architecture we defer to the original paper [9].

FullSubNet is composed of two models that operate on different parts of the noisy magnitude spectrogram $\mathbf{X} \in \mathbb{R}^{T \times F}$, which is the real-valued output of the short-time Fourier transform over $\mathbf{x}$. The model trains to minimize the $L_1$ or $L_2$ loss between complex Ratio Mask $\mathbf{M} \in \mathbb{R}^{T \times F}$ and the target complex Ideal Ratio Mask [35]. The mask is multiplied elementwise with the noisy spectrogram to obtain enhanced spectrogram $\hat{\mathbf{Y}} = \mathbf{M} \cdot \mathbf{X}$. Finally, inverse-STFT (iSTFT) gives us the desired waveform $\hat{\mathbf{y}}$.

For each model, we supplement the existing loss with our proposed eGeMAPS loss. Our estimator produces the eGeMAPS estimates from the enhanced spectrogram:

$$\hat{\mathbf{e}} = h_\phi(\hat{\mathbf{Y}}), \ \hat{\mathbf{e}} \in \mathbb{R}^D \quad (1)$$

where $h_\phi$ is our estimator parameterized by weights $\phi$. We apply STFT on $\hat{\mathbf{y}}$ to get $\hat{\mathbf{Y}}$, which allows us to apply our method to any enhancement model.

We also define the original non-differentiable function to calculate eGeMAPS given by [30], which is implemented in the openSMILE[34] toolkit:

$$\mathbf{e} = o(\mathbf{y}), \mathbf{e} \in \mathbb{R}^D \quad (2)$$

The toolkit takes in waveform $\mathbf{y}$ rather than spectrogram $\mathbf{Y}$. Given the difficulty of modeling waveforms, our estimator operates on spectrograms. Our training results for estimating eGeMAPS from spectrograms validate this decision.

For each batch of size $B$, our loss is the mean squared error between estimated eGeMAPS from enhanced speech and ground-truth eGeMAPS from clean speech:

$$\mathcal{L}_{\theta,\phi}^{\text{eGeMAPS}} = \frac{1}{B} \sum_{i=1}^{B} \|\hat{\mathbf{e}}_i^{\text{enh}} - \mathbf{e_i}^{\text{clean}}\|_2^2 \quad (3)$$

where superscripts enh and clean denote enhanced and clean speech, respectively. We combine each model's original loss with our loss:

$$\mathcal{L}_{\theta,\phi}^{\text{final}} = \mathcal{L}_\theta^{\text{original}} + \lambda \mathcal{L}_{\theta,\phi}^{\text{eGeMAPS}} \quad (4)$$

Note that $\mathcal{L}^{\text{eGeMAPS}}$ depends on both the enhancement model parameters $\theta$ *and* estimator parameters $\phi$, while the original SE model objectives only depend on $\theta$. Based on $\mathcal{L}_{\theta,\phi}^{\text{final}}$ we fine-tune the enhancement weights $\theta$ to optimize the original and eGeMAPS objectives. We also experiment with fixing and fine-tuning $\phi$, which we discuss further in Section 4.4. Figure 2 visualizes this fine-tuning scheme for arbitrary enhancement models.

## 4. Experiments

### 4.1. Data

We used Deep Noise Suppression (DNS) Challenge 2020 [32] data synthesized from their provided script, which creates 12,000 utterances of 30 seconds (s). We synthesized another 1,200 utterances for a validation set. The clean samples are from Librivox[3], and added noise clips are from Audioset[4] and Freesound[5]. The Signal to Noise Ratio (SNR) is sampled uniformly between 0 and 40 decibels (dB). Demucs and FullSubnet both train on sub-samples of the 30s audios. We followed

---
[3] https://librivox.org/
[4] https://research.google.com/audioset/
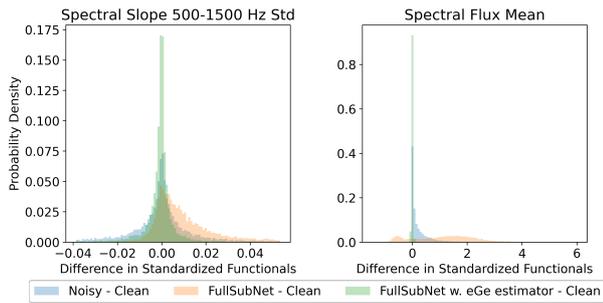[5] https://freesound.org/

Figure 4: *Test data distribution of eGeMAPS differences after fine-tuning, showing that we have pushed the differences towards 0. In the right-hand figure, the fine-tuning corrects the feature that was worse than noisy speech.*

| Model | PESQ | STOI | DNSMOS |
|---|---|---|---|
| Clean | - | - | 4.02 |
| Noisy | 1.58 | 91.52 | 3.16 |
| FullSubNet [31] | 2.89 | 96.41 | 3.91 |
| FullSubNet + eGeMAPS estimator | **2.91** | **96.43** | **3.93** |
| Demucs [9] | 2.65 | 96.54 | **3.80** |
| Demucs + eGeMAPS estimator | **2.81** | **96.83** | 3.78 |

Table 1: *Experimental results showing the benefit of eGeMAPS loss on FullSubNet and Demucs, measured by standard perceptual metrics. PESQ and STOI for clean are blank because the metrics are calculated with clean as reference.*

Demucs' configuration: for each 30 second audio, we took four seconds of audio with one second stride to create 27 samples. We used the DNS synthetic test set with no reverberation. This set consists of 150 utterances from Graz University's clean speech dataset [36], combined with noise categories chosen to represent real-world noise interruptions to audio calls. The SNR was sampled uniformly between 0 and 25 dB.

### 4.2. Estimator Specifications

We obtained our eGeMAPS estimator by first pre-training a Predictive Aux-VAE [37] on the clean spectrograms of our dataset. Pre-training provides quality representations of spectrograms and allows for easier optimization of downstream prediction of eGeMAPS. We removed the original VAE decoder after spectrogram pre-training. For our final eGeMAPS estimator, we used the VAE encoder followed by self-attention pooling [38] and two linear layers. The encoder is composed of four layers of convolutions, with the first three followed by tanh non-linearity and batch-normalization [39]. We confirm the estimator is performant by showing MSE values of $1.9 \cdot 10^{-4}$, $1.6 \cdot 10^{-4}$, and $2.6 \cdot 10^{-4}$ on train, validation, and test. Code for the estimator and fine-tuning of enhancement will be available at `https://github.com/muqiaoy/eGeMAPS_estimator`.

### 4.3. Results

In Figure 3 we see the success of our fine-tuning. Recall that eGeMAPS functionals are statistics of LLDs, such as mean and std. We show percent improvement in the means, where percent improvement is the reduction in difference from clean

| Setting | PESQ | STOI |
|---|---|---|
| $\lambda = 0$ | 2.65 | 96.54 |
| $\lambda = 0.1$ | 2.78 | 96.77 |
| $\lambda = 1$ | **2.81** | **96.83** |
| $\lambda = 10$ | 2.80 | 96.75 |
| $\lambda = 1$, fine-tuning $\phi$ | 2.79 | 96.75 |

Table 2: *An ablation study of the effect of eGeMAPS estimator on Demucs, where $\phi$ refers to the estimator weights as introduced in Section 3.3. If not specified, $\phi$ is fixed during the enhancement model fine-tuning stage.*

speech. Therefore 100% improvement means that the difference has been removed, and the functional matches clean. We see improvement in all features.

Figure 4 provides another way to visualize this improvement, through plotting differences between eGeMAPS as in Section 3.2. The green distribution on the right-hand side of the plot now shows close to ideal distribution of differences, completely peaked around 0. The left-hand side does not fully minimize the differences but shows good improvements compared to the original enhancement.

In Table 1 we show the effect of improvements in eGeMAPS differences on existing perceptual metrics. Note that PESQ and STOI are calculated with clean audios as reference, while DNSMOS [19] is directly evaluated on the enhanced audios. We report DNSMOS on clean to show our upper bound of performance. We see modest improvements in all metrics for FullSubNet. We also show large improvements in PESQ and STOI for Demucs. These indicate that our eGeMAPS estimator is improving the overall speech quality by the fine-grained speech characteristics.

### 4.4. Ablation Study of eGeMAPS Estimator

We analyzed multiple configurations for pairing the estimator with enhancement models, shown in Table 2. In all configurations, the original VAE encoder weights were fixed. We experimented with pre-training the final linear layers before fine-tuning the speech-enhancement model. We also tried training them while fine-tuning the enhancement model, but this did not converge. We found that fine-tuning the estimator along with the enhancement model performed worse in initial experiments, and therefore decided to perform our hyperparameter search with fixed estimator. We hypothesized that fine-tuning the estimator could add robustness to enhanced speech as input, but we conjecture that it creates a harder optimization problem.

We experimented with the weight multiplier, $\lambda$, of our eGeMAPS loss to match the scale of loss values, but found that $\lambda = 1$ is optimal for fine-tuning $\theta$ when using fixed $\phi$.

## 5. Conclusion

We demonstrate that fine-grained speech characteristics can significantly improve speech enhancement. We create a differentiable eGeMAPS estimator, allowing us to fine-tune existing models by minimizing acoustic parameter differences. The resulting enhanced audio yields superior results over traditional methods that marginally improve acoustic features and sometimes make them worse. Most perceptual evaluation metrics and acoustic functional statistics can be improved using our approach.

# 6. References

[1] S. Boll, "Suppression of acoustic noise in speech using spectral subtraction," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 27, no. 2, pp. 113–120, 1979.

[2] P. Scalart and J. Filho, "Speech enhancement based on a priori signal to noise estimation," in *1996 IEEE International Conference on Acoustics, Speech, and Signal Processing Conference Proceedings*, vol. 2, 1996, pp. 629–632 vol. 2.

[3] X. Lu, Y. Tsao, S. Matsuda, and C. Hori, "Speech enhancement based on deep denoising autoencoder," in *INTERSPEECH*, 2013.

[4] Y. Xu, J. Du, L.-R. Dai, and C.-H. Lee, "A regression approach to speech enhancement based on deep neural networks," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 23, no. 1, pp. 7–19, 2015.

[5] F. Weninger, F. Eyben, and B. Schuller, "Single-channel speech separation with memory-enhanced recurrent neural networks," in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2014, pp. 3709–3713.

[6] F. Weninger, H. Erdogan, S. Watanabe, E. Vincent, J. L. Roux, J. R. Hershey, and B. Schuller, "Speech enhancement with LSTM recurrent neural networks and its application to noise-robust ASR," in *LVA/ICA*, 2015.

[7] H. Zhao, S. Zarar, I. Tashev, and C. Lee, "Convolutional-recurrent neural networks for speech enhancement," *CoRR*, vol. abs/1805.00579, 2018. [Online]. Available: http://arxiv.org/abs/1805.00579

[8] S. Pascual, A. Bonafonte, and J. Serrà, "SEGAN: speech enhancement generative adversarial network," *CoRR*, vol. abs/1703.09452, 2017. [Online]. Available: http://arxiv.org/abs/1703.09452

[9] A. Defossez, G. Synnaeve, and Y. Adi, "Real time speech enhancement in the waveform domain," *arXiv preprint arXiv:2006.12847*, 2020.

[10] Y. Xu, J. Du, L.-R. Dai, and C.-H. Lee, "An experimental study on speech enhancement based on deep neural networks," *IEEE Signal Processing Letters*, vol. 21, no. 1, pp. 65–68, 2014.

[11] C. K. A. Reddy, E. Beyrami, J. Pool, R. Cutler, S. Srinivasan, and J. Gehrke, "A scalable noisy speech dataset and online subjective test framework," 2019. [Online]. Available: https://arxiv.org/abs/1909.08050

[12] J. Turian and M. Henry, "I'm sorry for your loss: Spectrally-based audio distances are bad at pitch," 2020. [Online]. Available: https://arxiv.org/abs/2012.04572

[13] P. Plantinga, D. Bagchi, and E. Fosler-Lussier, "Perceptual loss with recognition model for single-channel enhancement and robust ASR," *CoRR*, vol. abs/2112.06068, 2021. [Online]. Available: https://arxiv.org/abs/2112.06068

[14] S.-W. Fu, C.-F. Liao, Y. Tsao, and S.-D. Lin, "Metricgan: Generative adversarial networks based black-box metric scores optimization for speech enhancement," in *International Conference on Machine Learning*. PMLR, 2019, pp. 2031–2041.

[15] Y. Koizumi, K. Niwa, Y. Hioka, K. Kobayashi, and Y. Haneda, "Dnn-based source enhancement self-optimized by reinforcement learning using sound quality measurements," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 81–85.

[16] J. M. Martin-Doñas, A. M. Gomez, J. A. Gonzalez, and A. M. Peinado, "A deep learning loss function based on the perceptual evaluation of the speech quality," *IEEE Signal Processing Letters*, vol. 25, no. 11, pp. 1680–1684, 2018.

[17] A. Rix, J. Beerends, M. Hollier, and A. Hekstra, "Perceptual evaluation of speech quality (pesq)-a new method for speech quality assessment of telephone networks and codecs," in *2001 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No.01CH37221)*, vol. 2, 2001, pp. 749–752 vol.2.

[18] C. H. Taal, R. C. Hendriks, R. Heusdens, and J. Jensen, "An algorithm for intelligibility prediction of time–frequency weighted noisy speech," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 7, pp. 2125–2136, 2011.

[19] C. K. Reddy, V. Gopal, and R. Cutler, "Dnsmos: A non-intrusive perceptual objective speech quality metric to evaluate noise suppressors," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 6493–6497.

[20] T. Manjunath, "Limitations of perceptual evaluation of speech quality on voip systems," in *2009 IEEE International Symposium on Broadband Multimedia Systems and Broadcasting*, 2009, pp. 1–6.

[21] H. Zhang, X. Zhang, and G. Gao, "Training supervised speech separation system to improve stoi and pesq directly," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 5374–5378.

[22] Y. Koizumi, K. Niwa, Y. Hioka, K. Kobayashi, and Y. Haneda, "Dnn-based source enhancement to increase objective sound quality assessment score," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 26, no. 10, pp. 1780–1792, 2018.

[23] M. R. Saddler, A. Francl, J. Feather, K. Qian, Y. Zhang, and J. H. McDermott, "Speech denoising with auditory models," *arXiv preprint arXiv:2011.10706*, 2020.

[24] T.-A. Hsieh, C. Yu, S.-W. Fu, X. Lu, and Y. Tsao, "Improving Perceptual Quality by Phone-Fortified Perceptual Loss Using Wasserstein Distance for Speech Enhancement," in *Proc. Interspeech 2021*, 2021, pp. 196–200.

[25] S. Kataria, J. Villalba, and N. Dehak, "Perceptual loss based speech denoising with an ensemble of audio pattern recognition and self-supervised models," in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 7118–7122.

[26] J. Hillenbrand, R. Cleveland, and R. Erickson, "Acoustic correlates of breathy vocal quality," *Journal of speech and hearing research*, vol. 37, pp. 769–78, 09 1994.

[27] G. d. Krom, "Some spectral correlates of pathological breathy and rough voice quality for different types of vowel fragments," *Journal of Speech, Language, and Hearing Research*, vol. 38, no. 4, pp. 794–811, 1995.

[28] H. Kasuya, S. Ogawa, Y. Kikuchi, and S. Ebihara, "An acoustic analysis of pathological voice and its application to the evaluation of laryngeal pathology," *Speech Communication*. [Online]. Available: https://www.sciencedirect.com/science/article/pii/0167639386900063

[29] L. F. Gallardo and B. Weiss, "Perceived interpersonal speaker attributes and their acoustic features," 2017.

[30] F. Eyben, K. R. Scherer, B. W. Schuller, J. Sundberg, E. André, C. Busso, L. Y. Devillers, J. Epps, P. Laukka, S. S. Narayanan, and K. P. Truong, "The geneva minimalistic acoustic parameter set (gemaps) for voice research and affective computing," *IEEE Transactions on Affective Computing*, vol. 7, no. 2, pp. 190–202, 2016.

[31] X. Hao, X. Su, R. Horaud, and X. Li, "Fullsubnet: a full-band and sub-band fusion model for real-time single-channel speech enhancement," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 6633–6637.

[32] C. K. Reddy, V. Gopal, R. Cutler, E. Beyrami, R. Cheng, H. Dubey, S. Matusevych, 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.

[33] S.-W. Fu, C. Yu, T.-A. Hsieh, P. Plantinga, M. Ravanelli, X. Lu, and Y. Tsao, "Metricgan+: An improved version of metricgan for speech enhancement," 2021.

[34] F. Eyben, M. Wöllmer, and B. Schuller, "Opensmile: the munich versatile and fast open-source audio feature extractor," in *Proceedings of the 18th ACM international conference on Multimedia*, 2010, pp. 1459–1462.

[35] D. S. Williamson, Y. Wang, and D. Wang, "Complex ratio masking for monaural speech separation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 24, no. 3, pp. 483–492, 2016.

[36] G. Pirker, M. Wohlmayr, S. Petrik, and F. Pernkopf, "A pitch tracking corpus with evaluation on multipitch tracking scenario." 01 2011, pp. 1509–1512.

[37] S. Springenberg, E. Lakomkin, C. Weber, and S. Wermter, "Predictive Auxiliary Variational Autoencoder for Representation Learning of Global Speech Characteristics," in *Proc. Interspeech 2019*, 2019, pp. 934–938.

[38] P. Safari, M. India, and J. Hernando, "Self-attention encoding and pooling for speaker recognition," 2020. [Online]. Available: https://arxiv.org/abs/2008.01077

[39] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International conference on machine learning*. PMLR, 2015, pp. 448–456.