# DNSMOS: A NON-INTRUSIVE PERCEPTUAL OBJECTIVE SPEECH QUALITY METRIC TO EVALUATE NOISE SUPPRESSORS

*Chandan K A Reddy, Vishak Gopal, Ross Cutler*

Microsoft Corporation, Redmond, WA
chkarada@microsoft.com, vishak.gopal@microsoft.com, rcutler@microsoft.com

## ABSTRACT

Human subjective evaluation is the "gold standard" to evaluate speech quality optimized for human perception. Perceptual objective metrics serve as a proxy for subjective scores. The conventional and widely used metrics require a reference clean speech signal, which is unavailable in real recordings. Previous no-reference approaches correlate poorly with human ratings and are not widely adopted in the research community. One of the biggest use cases of these perceptual objective metrics is to evaluate noise suppression algorithms. This paper introduces a multi-stage self-teaching based perceptual objective metric that is designed to evaluate noise suppressors. The proposed method generalizes well in challenging test conditions with a high correlation to human ratings.

***Index Terms***— Speech, Perceptual Speech Quality, Objective Metric, Deep Noise Suppressor, Metric.

## 1. INTRODUCTION

Subjective evaluation of speech quality is the most reliable way to evaluate Speech Enhancement (SE) methods [1]. However, subjective tests are not scalable as they require a considerable number of listeners, the process is laborious, time-consuming, and expensive. Conventional objective speech quality metrics such as Perceptual Evaluation of Speech Quality (PESQ) [2], Perceptual Objective Listening Quality Analysis (POLQA) [3] and Signal to Distortion Ratio (SDR) are widely used to evaluate Speech Enhancement (SE) algorithms optimized for human perception. Some of these metrics are designed to predict the subjective Mean Opinion Score (MOS) obtained using the Telecommunication Standardization Sector of the International Telecommunication Union (ITU-T) Recommendation P.800 [4]. There are also the intrusive metrics requiring a reference clean speech to compute the quality score. However, they are shown to correlate poorly with human rating when used for SE tasks that involve perceptually invariant transformations [1]. Also, intrusive metrics cannot be used to evaluate real recordings when clean reference is unavailable in realistic scenarios.

ITU-T Recommendation P.563 is a non-intrusive technique and can directly operate on the degraded signal [5].

However, it was developed for narrow-band applications and works well on limited impairment types. Recently, Deep Neural Networks (DNN) based approaches have been proposed to estimate the speech quality scores [6, 7, 8]. Some of these learning-based approaches use other objective metrics as the ground truth to train their speech quality predictor. A handful of published methods use MOS obtained using P.800 as the ground truth to train their models. In [9], the authors trained the model to identify the Just Noticeable Difference (JND). MOS predictors trained on actual human ratings are more reliable than the ones trained to predict other objective metrics like PESQ or POLQA. The accuracy and robustness of the learned models depend on the quality of the human labels and also the quantity and diversity of the audio clips. There is no large scale human-labeled data set publicly available to train a model that can robustly generalize to a variety of audio impairments, especially for SE task. Though the learning-based approaches have a higher correlation to human ratings when tested on a limited and small test set, none of these speech quality metrics are as widely used as some of the intrusive objective metrics. Most of the publications in the area of noise suppression and SE test their methods on synthetic noisy speech test sets using intrusive speech quality metrics such as PESQ, POLQA and SDR.
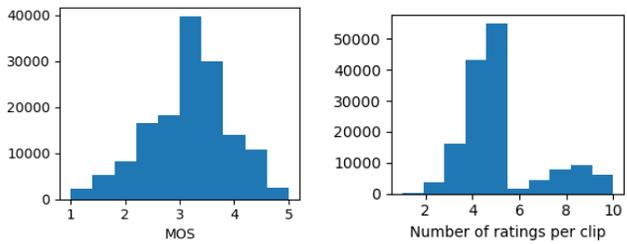
In this paper, we have developed a robust objective perceptual speech quality metric called DNSMOS. It can be used to stack rank different Deep Noise Suppression (DNS) methods based on MOS estimates with great accuracy and hence the name DNSMOS. A simple Convolutional Neural Network (CNN) based model is trained using the ground truth human ratings obtained using ITU-T P.808 [10, 11]. P.808 is an online subjective testing framework that is highly reproducible and shown to stack rank noise suppression models with high accuracy when each model is tested as an average over a statistically significant number of clips. However, P.808 ratings per clip can be noisy due to various factors. The absolute MOS per clip might vary for the same clip in different P.808 runs due to the inherent noise in human rating [12]. The number of ratings per clip varies due to spam removal. Also, all the spammers might not be filtered due to the limitations of spam filtering logic in P.808 implementation. We train a multi-stage self-teaching model inspired by *continual*

*life long learning* [13] to learn in the presence of noisy labels. We show that the multi-stage training significantly improves the correlation and generalizes remarkably well to other impairment types that are very different from what was used in the training set. We found the DNSMOS very useful in our audio/speech research as it generalizes well to a variety of speech impairments. Hence, we are providing the DNSMOS as an Azure service for other researchers to use. The details of the API are at `www.microsoft.com/en-us/research/dns-challenge/dnsmos`.
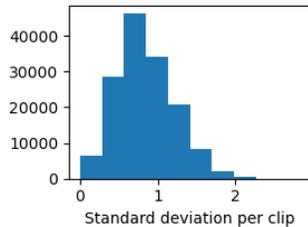
## 2. DATA AND SUBJECTIVE RATINGS

We used 600 noisy speech test clips comprised of a combination of synthetic and real recordings. The synthetic test set is a combination of both non-reverberant and reverberant clips. The real recordings were captured in a variety of noise types and Signal to Noise Ratio (SNR) and target levels. The test set is comprised of over 100 noise types and speakers. More details about the creation of these test sets can be found in [11]. These 600 clips were processed by over 200 different noise suppression methods. Some of these methods improved the overall speech quality, but they also introduced a variety of artifacts as a consequence of over or under suppression of noise, resulting in speech and noise distortions. The speech quality ratings of the processed clips varied from very poor (MOS=1) to excellent (MOS=5). The distribution of the MOS scores in the training data is shown in Figure 1a. The scores are highly skewed with most ratings populated in the range 3<MOS<4 and fewer ratings in both the tails. The distributions in Figure 1a also shows that the current state of the art noise suppression techniques are not good enough to shift the overall distribution above MOS of 4. Figure 1b shows the distribution of the number of ratings per clip in the training set. The skewness in the distribution is due to the fact that some of the ratings were removed in the process of removing spam ratings. Also, we used 10 ratings per clip in some of the experiments and 5 ratings per clip in others. Figure 1c shows the distribution of standard deviation per clip. There are a significant number of clips with a standard deviation of over 1 MOS, which shows that the audio clips are highly subjective in nature. Increasing the number of ratings per clip will reduce the standard deviation, but also increases the cost.

The subjective human ratings are obtained in several P.808 runs conducted over several months. Multiple noise suppression methods are compared in each P.808 run. Each P.808 run included the best performing noise suppressor, original noisy speech, and a couple of methods with intermediate perceptual quality from previous runs as anchors. Hence, some of the clips were rated multiple times. In total, we have about 120,000 audio clips with associated MOS scores as ground truth. The average length of each audio clip was about 9 seconds. In [10], we show that P.808 is highly reproducible with high SRCC between runs when we average the ratings across



(a) Distribution of MOS ratings    (b) Distribution of number of votes



(c) Distribution of standard deviation per clip

**Fig. 1**: Statistics of the training set

clips processed by a given model. However, the absolute values of the scores vary between P.808 runs as the raters differ and humans tend to be inconsistent in perceptual tasks due to various biases.

## 3. DNSMOS

### 3.1. Learning in the presence of label biases

The variation in the absolute values between P.808 runs might pose a challenge in training models that learn to predict MOS scores when the training data includes the same audio clip with different ratings. We know from the literature that DNN models can generalize well in spite of the noise for classification tasks [14, 15]. The variation in the absolute scores from different runs can be treated as noise in the labels. Also, the clips with fewer ratings can be treated as noise due to higher standard deviation. There are other unknown human biases that play a role in the variation of ratings. We initially took a direct approach where we trained a deep model by including all the available data without any modification. The model tends to learn the weighted average of the ratings across runs and also tries to average out other biases. This plain approach gave a good correlation with human ratings when tested on the data that was very similar to the training set. However, the correlation dropped when tested on a more challenging test set that was very different from the training data.

We improved the generalization and accuracy of the predictor by training a few stages of self-teaching models. In [16], a sequence of self-teachers were used to improve the

generalizability of a classification task with binary cross-entropy loss in noisy and weakly labeled conditions. We use a similar approach for the regression problem at hand with Mean Squared Error (MSE) as the loss function. Unlike knowledge distillation where a larger teacher model is used to produce soft labels to train a smaller student model, we use the same model architecture for both teacher and the student, and hence the name self-teacher. The student model uses the weighted average of the predictions from one or multiple teacher models and the original human MOS ratings. The primary teacher model $M_0$ is trained until the loss saturates using the original ground truth human MOS ratings $r$. The student model $M_s$ at stage $s$, is trained using the new target given by,

$$\hat{r}_s = \alpha_0 r + \sum_{i=0}^{s-1} \alpha_{i+1} \hat{r}_i \qquad (1)$$

where $\hat{r}_i$ is the prediction of model $M_i$ and $\sum_{i=0}^{s} \alpha_i = 1$. In this work, we restrict $s$ to 2. With some assumptions to the noise distribution, it can be theoretically proven that the $M_s$ is better than $M_0$ for $s > 0$. Since the distribution of noise in our case is complex, we show through empirical experiments that the predicted MOS of $M_i$ is of higher accuracy than $M_0$ for appropriately chosen values of $\alpha$.

## 3.2. Features

Recently, researchers have seen success in learning features within the model for tasks such as SE [17], speech, and music synthesis [18] and to learn acoustic models [19]. They show that using the time-domain waveform requires a larger model trained on a larger and diverse data set to ensure generalization. The ground truth MOS scores are obtained for audio clips with an average length of 9 secs sampled at 16 kHz. This leads to a very large input dimension if we are treating it as a vector and the model requires many layers to compress and extract input features. We used Log power Mel spectrogram as input feature as it correlates well with human perception and is proven to work very well for analyzing speech quality [7]. For spectral features, we used a frame size of 20 ms with a hop length of 10 ms and 120 Mel frequency bands. The input features are then converted to dB scale.

## 3.3. Prediction model

For predicting the MOS scores, we explored different configurations of CNN based models. The architecture for the best performing model is shown in Table 1, which is the architecture for all $M$. The input to the model is log power Mel Spectrogram with 120 Mel bands computed over a clip of length 9 secs sampled at 16 kHz with a frame size of 20 ms and hop length of 10 ms. This results in an input dimension of 900 x 120. The model was trained with a batch size of 32 using the Adam optimizer and MSE loss function until the loss saturated. We experimented by adding batch normalization

**Table 1**: DNSMOS Prediction Model

| Layer | Output dimension |
|---|---|
| Input | 900 x 120 x 1 |
| Conv: 32, (3 x 3), 'ReLU' | 900 x 120 x 32 |
| MaxPool: (2 x 2), Dropout(0.3) | 450 x 60 x 32 |
| Conv: 32, (3 x 3), 'ReLU' | 450 x 60 x 32 |
| MaxPool: (2 x 2), Dropout(0.3) | 225 x 30 x 32 |
| Conv: 32, (3 x 3), 'ReLU' | 225 x 30 x 32 |
| MaxPool: (2 x 2), Dropout(0.3) | 112 x 15 x 32 |
| Conv: 64, (3 x 3), 'ReLU' | 112 x 15 x 64 |
| GlobalMaxPool | 1 x 64 |
| Dense: 64, 'ReLU' | 1 x 64 |
| Dense: 64, 'ReLU' | 1 x 64 |
| Dense: 1 | 1 x 1 |

layers after every Conv layer in Table 1. However, adding batch normalization reduces the prediction accuracy of low volume clips. Humans tend to give lower ratings to the clips with low amplitudes. We want the model to capture the variations in the target levels of the data. Hence, we avoid any kind of feature normalization. We also explored different network architectures including CNN followed by LSTM. The model in Table 1 generalized the best and was of least complexity.
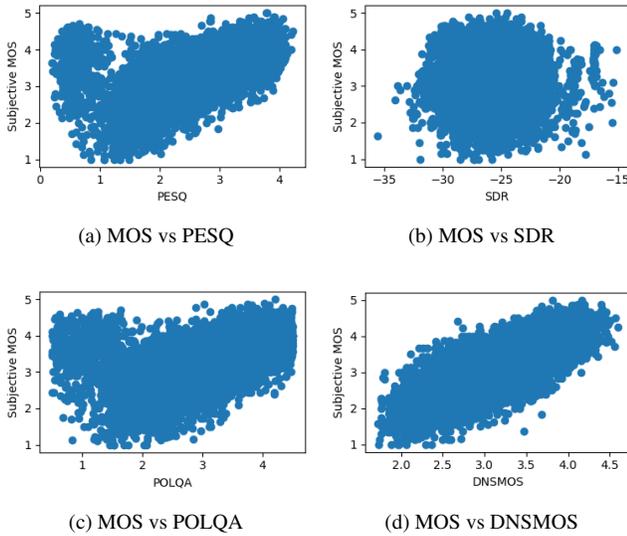
## 4. EXPERIMENTAL RESULTS

### 4.1. Evaluation metric

Pearson Correlation Coefficient (PCC) or MSE between the predictions of the developed objective metric and the ground truth human ratings is commonly used to measure the accuracy of the model [7, 8]. From the earlier discussion, we know that P.808 correlation is highly repeatable between runs when averaged across a set of clips per condition, which can be formed by grouping clips enhanced by a particular SE model or based on other criteria like SNR or reverb RT60 times. The PCC computed on the average of ratings per group across different runs is >0.9. We also found that PCC computed on the same clips but from two different P.808 runs is only about 0.5 due to the high rating noise per clip.
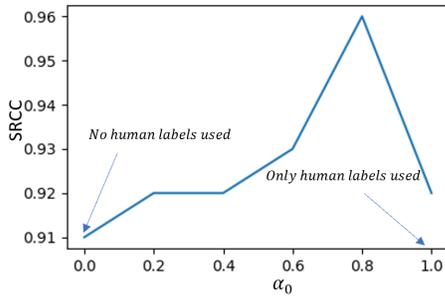
Hence, for stack ranking different noise suppressors we evaluate by computing the average of ratings across the entire test set for each model. Therefore, we compute SRCC and

**Table 2**: Correlation of DNSMOS with other widely used objective metrics

| | PESQ | SDR | POLQA | DNSMOS ($M_0$) |
|---|---|---|---|---|
| PCC | 0.78 | 0.23 | 0.79 | 0.93 |
| SRCC | 0.82 | 0.25 | 0.84 | 0.94 |

(a) MOS vs PESQ　　　(b) MOS vs SDR

(c) MOS vs POLQA　　　(d) MOS vs DNSMOS

**Fig. 2**: Scatter plots of different objective metrics vs MOS



**Fig. 3**: SRCC between MOS and $M_1$ for different $\alpha_0$ values

PCC between averaged human ratings and averaged DNSMOS per model. SRCC gives us the stack ranking accuracy of various SE models.

## 4.2. Results

### 4.2.1. DNSMOS vs other objective speech quality metrics

Table 2 shows the PCC and SRCC between human ratings and widely used PESQ, POLQA, SDR and DNSMOS ($M_0$) computed on the Interspeech DNS challenge blind test results [11]. The results show that SDR correlates poorly with MOS. PESQ and POLQA reasonable PCC and SRCC. But the primary DNSMOS model ($M_0$) correlates significantly better than other objective metrics and is highly reliable in stack ranking the SE models. The scatter plots in Figure. 2 shows that DNSMOS aligns with MOS better than SDR, PESQ, and POLQA. Both PESQ and POLQA have a cluster on the left top of the graph indicating that they tend to penalize certain artifacts more than we humans do.

**Table 3**: Per category SRCC between ($M_0$) and MOS on DNS Challenge 2 results

| $M_i(\alpha_0,\alpha_1,\alpha_2)$ | English | Non-English | Tonal | Emotional | Singing | Overall |
|---|---|---|---|---|---|---|
| $M_0(1,-,-)$ | 0.78 | 0.67 | 0.51 | 0.69 | 0.36 | 0.65 |
| $M_1(0.8, 0.2, -)$ | 0.81 | 0.86 | 0.81 | 0.58 | 0.61 | **0.90** |
| $M_1(0.6, 0.4, -)$ | 0.90 | 0.85 | 0.82 | 0.72 | 0.58 | 0.88 |
| $M_1(0.4, 0.6, -)$ | **0.93** | **0.92** | **0.87** | 0.67 | 0.57 | 0.88 |
| $M_1(0.2, 0.8, -)$ | 0.90 | 0.85 | 0.82 | 0.72 | 0.58 | 0.88 |
| $M_2(-, 0.5, 0.5)$ | 0.90 | 0.86 | 0.81 | 0.75 | 0.54 | **0.90** |
| $M_2(0.6, 0.2, 0.2)$ | 0.91 | 0.85 | 0.82 | **0.78** | **0.72** | 0.89 |
| $M_2(0.8, 0.1, 0.1)$ | 0.88 | 0.85 | 0.82 | 0.68 | 0.69 | 0.88 |

### 4.2.2. Generalizability of DNSMOS

We initially analyzed the effect of $\alpha_0$ on $M_1$ to understand the weightage given to the teacher model $M_0$ vs human ratings. Figure. 3 shows the SRCC between MOS and $M_1$ for different values of $\alpha_0$ computed on one of the P.808 runs with 12 closely stack ranked models. We see in Figure. 3 that the SRCC is higher for $\alpha_0 = 0.8$.

The noisy and the processed clips comprised of a similar distribution as that of the training set. In order to test the generalization capability of training multiple stages of self-teaching models, we evaluated on the ICASSP 2021 DNS Challenge 2 final submissions [20], which comprised of the blind test set processed by 18 different noise suppressors. The blind set included emotional speech, singing, utterances in English and non-English languages including a few tonal languages. The DNSMOS models were trained on the data which was predominantly in English with very little or no emotional content such as laughter, crying, yelling, anger, etc. The training set did not include singing and other non-English languages as well. Table 3 shows the SRCC between the DNSMOS models at various stages and P.808 MOS for each category in the blind set. The SRCC of the primary model $M_0$ is quite low as it does not generalize well for categories other than English. $M_1$ significantly improved SRCC from 0.65 to 0.9 for $\alpha_0 = 0.8$. The optimal values of $\alpha_0$ were different for each category in table 3. We also experimented with $M_2$ for a couple of settings of $(\alpha_0,\alpha_1,\alpha_2)$. For $M_2(-, 0.5, 0.5)$ with $M_1(0.2, 0.8, -)$ for stage 2 gave the best results. The other settings did not show much improvement over $M_1$.

## 5. CONCLUSION AND FUTURE WORK

DNSMOS is a robust speech quality metric designed to stack rank noise suppressors with great accuracy. The multi-stage self-teaching approach significantly improved accuracy. In the future, we would like to deepen our theoretical understanding of how the choice of $\alpha_i$ influences the accuracy of $M_i$ for different datasets. This will help in training DNSMOS on other impairment types such as network distortions, codec artifacts, and reverberation without the need for large scale data with human labels.

# 6. REFERENCES

[1] Chandan KA Reddy, Ebrahim Beyrami, Jamie Pool, Ross Cutler, Sriram Srinivasan, and Johannes Gehrke, "A scalable noisy speech dataset and online subjective test framework," *Proc. Interspeech 2019*, pp. 1816–1820, 2019.

[2] 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 *2001 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No.01CH37221)*, 2001, vol. 2, pp. 749–752 vol.2.

[3] John Beerends, Christian Schmidmer, Jens Berger, Matthias Obermann, Raphael Ullmann, Joachim Pomy, and Michael Keyhl, "Perceptual objective listening quality assessment (polqa), the third generation itu-t standard for end-to-end speech quality measurement part ii-perceptual model," *AES: Journal of the Audio Engineering Society*, vol. 61, pp. 385–402, 06 2013.

[4] "Itu-t recommendation p.800: Methods for subjective determination of transmission quality, itu-t recommendation p.800," Feb 1998.

[5] "Itu-t, recommendation p.563: Single-ended method for objective speech quality assessment in narrowband telephony applications," 2004.

[6] Xuan Dong and Donald S Williamson, "An attention enhanced multi-task model for objective speech assessment in real-world environments," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 911–915.

[7] Hannes Gamper, Chandan KA Reddy, Ross Cutler, Ivan J Tashev, and Johannes Gehrke, "Intrusive and non-intrusive perceptual speech quality assessment using a convolutional neural network," in *2019 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2019, pp. 85–89.

[8] A. R. Avila, H. Gamper, C. Reddy, R. Cutler, I. Tashev, and J. Gehrke, "Non-intrusive speech quality assessment using neural networks," in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 631–635.

[9] Pranay Manocha, Adam Finkelstein, Zeyu Jin, Nicholas J Bryan, Richard Zhang, and Gautham J Mysore, "A differentiable perceptual audio metric learned from just noticeable differences," *arXiv preprint arXiv:2001.04460*, 2020.

[10] Babak Naderi and Ross Cutler, "An open source implementation of itu-t recommendation p.808 with validation," 2020.

[11] Chandan KA Reddy, Vishak Gopal, Ross Cutler, Ebrahim Beyrami, Roger Cheng, Harishchandra Dubey, Sergiy Matusevych, Robert Aichner, Ashkan Aazami, Sebastian Braun, et al., "The interspeech 2020 deep noise suppression challenge: Datasets, subjective testing framework, and challenge results," *arXiv preprint arXiv:2005.13981*, 2020.

[12] ITUT Rec, "P. 1401, methods, metrics and procedures for statistical evaluation, qualification and comparison of objective quality prediction models," *International Telecommunication Union, Geneva, Switzerland*, 2012.

[13] German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter, "Continual lifelong learning with neural networks: A review," *Neural Networks*, vol. 113, pp. 54–71, 2019.

[14] Chandan KA Reddy, Ross Cutler, and Johannes Gehrke, "Supervised classifiers for audio impairments with noisy labels," *Proc. Interspeech 2019*, pp. 2573–2577, 2019.

[15] David Rolnick, Andreas Veit, Serge J. Belongie, and Nir Shavit, "Deep learning is robust to massive label noise," *ArXiv*, vol. abs/1705.10694, 2017.

[16] Anurag Kumar and Vamsi Krishna Ithapu, "A sequential self teaching approach for improving generalization in sound event recognition," .

[17] Alexandre Defossez, Gabriel Synnaeve, and Yossi Adi, "Real time speech enhancement in the waveform domain," 2020.

[18] Kou Tanaka, Hirokazu Kameoka, Takuhiro Kaneko, and Nobukatsu Hojo, "Wavecyclegan2: Time-domain neural post-filter for speech waveform generation," *arXiv preprint arXiv:1904.02892*, 2019.

[19] Y. Hoshen, R. J. Weiss, and K. W. Wilson, "Speech acoustic modeling from raw multichannel waveforms," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 4624–4628.

[20] Chandan K A Reddy, Harishchandra Dubey, Vishak Gopal, Ross Cutler, Sebastian Braun, Hannes Gamper, Robert Aichner, and Sriram Srinivasan, "Icassp 2021 deep noise suppression challenge," *arXiv preprint arXiv:2009.06122*, 2020.