

PITCH-AND-SPECTRUM-AWARE SINGING QUALITY ASSESSMENT WITH BIAS CORRECTION AND MODEL FUSION

Yu-Fei Shi, Yang Ai*, Ye-Xin Lu, Hui-Peng Du, Zhen-Hua Ling

National Engineering Research Center of Speech and Language Information Processing,
University of Science and Technology of China, Hefei, P. R. China

{zkddsr2023, yxlu0102, redmist}@mail.ustc.edu.cn, {yangai, zhling}@ustc.edu.cn

ABSTRACT

We participated in track 2 of the VoiceMOS Challenge 2024, which aimed to predict the mean opinion score (MOS) of singing samples. Our submission secured the first place among all participating teams, excluding the official baseline. In this paper, we further improve our submission and propose a novel Pitch-and-Spectrum-aware Singing Quality Assessment (PS-SQA) method. The PS-SQA is designed based on the self-supervised-learning (SSL) MOS predictor, incorporating singing pitch and spectral information, which are extracted using pitch histogram and non-quantized neural codec, respectively. Additionally, the PS-SQA introduces a bias correction strategy to address prediction biases caused by low-resource training samples, and employs model fusion technology to further enhance prediction accuracy. Experimental results confirm that our proposed PS-SQA significantly outperforms all competing systems across all system-level metrics, confirming its strong sing quality assessment capabilities.

Index Terms— sing quality assessment, MOS prediction, pitch histogram, bias correction, model fusion

1. INTRODUCTION

With the rapid development of singing voice synthesis (SVS) and singing voice conversion (SVC) systems, there is an urgent need for technology that can automatically assess the quality of generated singing voice, instead of traditional subjective listener scoring methods which are time-consuming and inefficient. However, in past research on singing quality assessment, most studies focus on the quality assessment of real recorded human singing voices [1, 2, 3]. To our knowledge, the quality assessment of generated singing voices has not yet been thoroughly investigated. Nonetheless, methods for speech quality assessment can serve as references and be applied to singing quality assessment. The mean opinion score (MOS) is the gold standard in the fields of speech synthesis and voice conversion [4], representing the average five-point rating given by humans to generated speech. Early MOS prediction models employed bidirectional long short-term memory recurrent neural network (BiLSTM-RNN) or convolutional neural network (CNN) to predict MOS score from input speech waveforms or amplitude spectra [5, 6, 7]. Recently, with the development of self-supervised learning (SSL) methods, fine-tuning SSL models and adapting them for speech MOS prediction has become one of the state-of-the-art approaches. MOS prediction is also applicable for singing quality assessment, but directly

transferring methods from speech MOS prediction is clearly inappropriate, as they may not align with the characteristics of singing voices.

The VoiceMOS Challenge 2024, launched this year, aims to encourage participants to build systems that predict MOS of singing voices generated by SVS and SVC systems in track 2. The challenge provides a platform for innovators to develop and test their models against a standardized dataset, advancing the progress of singing quality assessment. The organizers provide a singing quality evaluation dataset, SingMOS, along with a baseline. Participants are required to build systems to predict the MOS for the singing voices in the evaluation set, and the organizers rank the participating systems using certain metrics. Our submitted system secured first place among all participating teams (excluding the official baseline). Nevertheless, there is still significant room for improvement in the accuracy of singing MOS prediction.

To address the existing issues in singing quality assessment, this paper further improves our competition system submitted to track 2 of VoiceMOS Challenge 2024 and proposes a novel Pitch-and-Spectrum-aware Singing Quality Assessment (PS-SQA) method. To address the characteristics of singing voices, the PS-SQA innovatively introduces pitch-aware SSL-based MOS predictors and spectrum-aware SSL-based MOS predictors. They are built on the plain SSL-based MOS predictor by introducing pitch histograms and spectral-level acoustic features encoded by a non-quantized APCodec [8], respectively. These methods attempt to inject key singing information such as musical melody into the predictor, providing a new approach suitable for evaluating synthesized singing. These methods enable the predictor to have a more nuanced understanding of the melodic content of the singing, which is crucial for accurate quality assessment. Additionally, we noticed that the official SingMOS dataset has the issue of imbalanced training sample distribution. To address this, the PS-SQA introduces a bias correction branch, integrating into aforementioned predictors to mitigate the effects of such imbalances. Finally, PS-SQA also employs a model fusion strategy, comprehensively considering the results of multiple MOS predictors to provide a more thorough and accurate singing MOS score. Experimental results confirm that our proposed PS-SQA significantly improves our submission system and clearly outperforms all competing systems in track 2 of VoiceMOS Challenge 2024 in terms of the system-level spearman rank correlation coefficient (SRCC) used for ranking.

This paper is organized as follows: In Section 2, We provide a brief review of the related work involved in PS-SQA, encompassing SSL-based MOS predictors, pitch histograms, and advanced neural audio codecs. In Section 3, we give details of the construction and workflow of the various components that make up PS-SQA. In

*Corresponding author. This work was funded by the National Nature Science Foundation of China under Grant 62301521, the Anhui Provincial Natural Science Foundation under Grant 2308085QF200, and the Fundamental Research Funds for the Central Universities under Grant WK2100000033.

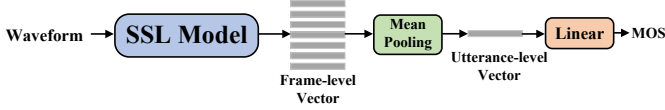


Fig. 1: A block diagram of a plain SSL-based MOS predictor.

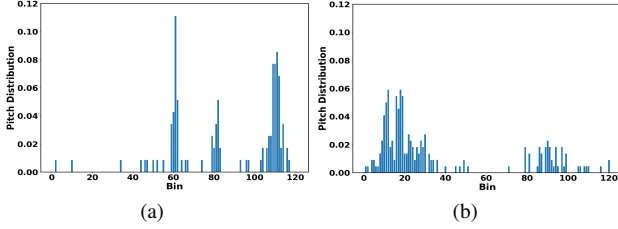


Fig. 2: The pitch histograms of (a) a good singing voice with MOS of 5.0 and (b) a poor singing voice with MOS of 2.4.

Section 4, we present our experimental results. Finally, we give conclusions in Section 5.

2. RELATED WORK

2.1. SSL-based MOS Predictor

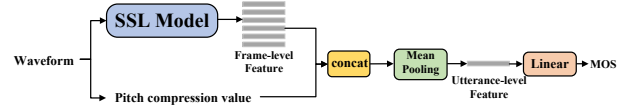
Recently, SSL models trained with a large amount of unlabeled data using self-supervised learning have been applied to MOS prediction, achieving impressive results [9]. In both the VoiceMOS Challenge 2022 and 2023, top-ranking teams employed fine-tuning on SSL models to achieve perfect predictive accuracy [10, 11]. Therefore, in the VoiceMOS Challenge 2024, we also fine-tuned the official baseline SSL-based MOS predictor to develop our system. The process of predicting MOS using SSL-based predictor is illustrated in Figure 1. The waveform is processed through a pre-trained SSL model to produce a frame-level feature vector, which is then averaged using mean-pooling to obtain an utterance-level one. Finally, a linear layer reduces the feature dimensionality to 1 to derive the corresponding MOS score. Assuming \hat{y} is the predicted MOS and y is its corresponding label, the loss function is defined as the L1 error between \hat{y} and y , i.e.,

$$\mathcal{L}_{ssl} = \|\hat{y} - y\|_1. \quad (1)$$

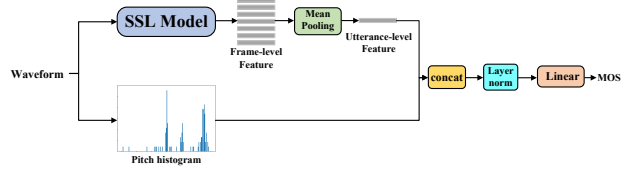
2.2. Pitch Histogram

Pitch is a critical measure in evaluating singing quality. Unlike direct use of detected fundamental frequency values in speech processing, assessing musical quality often involves transforming pitch values. Since adjacent notes in sheet music have consistent pitch ratios, folding identified pitch values over an octave can accurately reconstruct the melody of a singing voice. Specifically, in the MIDI scale, a complete octave contains 12 semitones, with the pitch ratio between adjacent semitones being $2^{1/12}$ and the pitch frequency ratio between each octave is 2. Additionally, we further divide the interval between adjacent semitones into 100 cents and the pitch frequency ratio between adjacent cents is $2^{1/1200}$. On this basis, we convert the aforementioned pitch frequency from the Hz scale (f_{Hz}) to the cent scale (f_{cent}) using the following formula:

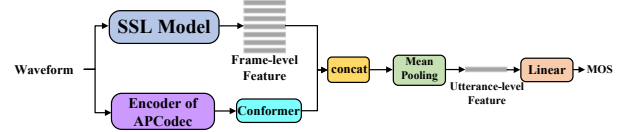
$$f_{cent} = 1200 \times \log_2 \frac{f_{Hz}}{440}, \quad (2)$$



(a) Compressed-pitch-aware SSL-based MOS predictor



(b) Pitch-histogram-aware SSL-based MOS predictor



(c) Spectrum-aware SSL-based MOS predictor

Fig. 3: Architectures of the pitch-aware and spectrum-aware SSL-based MOS predictors.

where 440 Hz (pitch-standard musical note A4) is considered as the base frequency.

In our implementation, we first use *PyWORLD* to extract the pitch in Hz (i.e., f_{Hz}) of different frames of singing. Then, using Equation 2, we get the converted pitch f_{cent} in the unit of cents. This allows us to obtain a pitch sequence that better reflects the characteristics of the singing voice. Following prior work [1, 12, 13], we treat 10 cents as the smallest counting unit (i.e., 1 bin). Therefore, there are a total of 12 semitones \times 10 bins = 120 bins in an octave. Then, we further convert the pitch value f_{cent} to a compressed continuous value $I(f_{cent}) \in [0, 120)$, i.e.,

$$I(f_{cent}) = \frac{f_{cent}}{10} \bmod 120, \quad (3)$$

where mod denotes the modulo operation.

In this way, pitch values from different frames of a singing voice can all be mapped into one octave. Subsequently, the pitch histogram $\mathbf{P} = \{P_1, \dots, P_j, \dots, P_{120}\}$ of a sample is obtained by calculating the ratio of the frame-level pitch count in each bin to the total number of frames, i.e.,

$$P_j = \frac{1}{N} \sum_{n=1}^N m_j^n, \quad (4)$$

where

$$m_j^n = \begin{cases} 1, & j-1 \leq I(f_{cent}^n) < j \\ 0, & \text{otherwise} \end{cases}. \quad (5)$$

Here, f_{cent}^n represents the pitch value in cent of the n -th frame and N is the total number of frames of a sing sample.

The melody of a song is typically composed of a set of pitch values that frequently appear in the vocals. In previous work [1, 14], it was found that for different singers performing the same song, the pitch histograms of good singers exhibit sharper peaks. This indicates that the notes of the song are consistently hit. In contrast, the pitch histograms of poor singers do not show such prominent peaks because they fail to consistently hit the dominant notes, resulting in being out of tune. Figure 2 shows pitch histograms extracted

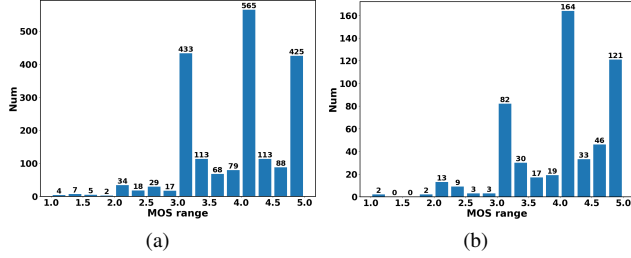


Fig. 4: Histograms of the number of samples in different MOS intervals for (a) training set and (b) validation set of SingMOS dataset.

from two samples provided in track 2 of VoiceMOS Challenge 2024, where it can be seen that the histogram of the high-MOS singing voice has more pronounced peaks compared to that of the low-MOS singing voice. This is why we chose the pitch histogram to assist in singing quality assessment in PS-SQA.

2.3. Neural Audio Codec

Neural audio codec is a crucial signal processing technology that compresses audio signals into discrete codes and then reconstructs the original audio from these codes by neural networks. It can also be considered a feature extractor, capable of effectively extracting intermediate representations that contain acoustic information from the audio, which can be used as input for MOS prediction models. Early neural audio codecs, such as SoundStream [15] and Encodec [16], directly encode the time-domain audio waveforms. Recently, Ai *et al.* proposes APCodec [8], a parametric neural audio codec that uses the audio amplitude and phase spectra as coding objects. Therefore, through explicit spectral modeling, the intermediate features encoded by APCodec contain richer spectral-level acoustic information. This is related to several amplitude-spectrum-based MOS prediction methods [7, 17], but it leverages the missing phase spectrum information in these methods, making it more suitable for MOS prediction. This is why we chose APCodec to extract spectral features to assist in singing quality assessment in PS-SQA.

3. PROPOSED METHOD

The core of PS-SQA lies in introducing pitch-aware SSL-based MOS predictors and spectrum-aware SSL-based MOS predictors based on the plain SSL-based MOS prediction framework, tailored for quality assessment that suits the characteristics of singing voices. Additionally, to overcome the issues caused by imbalanced training data, PS-SQA innovatively introduces a bias correction strategy. Finally, PS-SQA also employs model fusion techniques, aggregating the results of multiple top-ranking MOS predictors to output a comprehensive MOS score, further enhancing the predictive accuracy of PS-SQA. All predictors are trained using the loss function defined in Equation 1.

3.1. Pitch-aware SSL-based MOS predictor

As mentioned in Section 2.1, the SSL models are widely used in speech MOS prediction. The SSL models extract semantic information from speech waveforms for MOS prediction, which is clearly not suitable for singing quality assessment. Considering the strong relevance of pitch to singing quality discussed in Section 2.2, we propose to integrate pitch information into the plain SSL-based MOS predictor to enhance MOS prediction accuracy for singing quality

assessment. By explicitly providing pitch-related information to the network, this approach also releases network degrees of freedom to focus on learning non-pitch related properties.

Therefore, the PS-SQA designs pitch-aware SSL-based MOS predictors by incorporating pitch information into the plain SSL-based MOS prediction framework in a specific manner. Furthermore, by adopting different forms of pitch information, we attempt to construct a compressed-pitch-aware SSL-based MOS predictor and a pitch-histogram-aware SSL-based MOS predictor, respectively.

- **Compressed-pitch-aware SSL-based MOS predictor:** As shown in Figure 3(a), the compressed-pitch-aware SSL-based MOS predictor integrates the compressed value-constrained pitch sequence $\mathbf{f}_c = [I(f_{cent}^1), \dots, I(f_{cent}^N)]^\top$ into a plain SSL-based MOS prediction framework. Specifically, the SSL model extracts frame-level features from the waveform using the same frame shift as *PyWORLD* does when extracting the pitch from the waveform. This ensures that the two extracted sequences have the same temporal resolution. Then, we concatenate the SSL-model-processed frame-level features with \mathbf{f}_c along the dimension axis. After pooling along the time axis to obtain the utterance-level feature, we pass it through a linear layer to output the 1-dimensional MOS score.
- **Pitch-histogram-aware SSL-based MOS predictor:** As shown in Figure 3(b), the pitch-histogram-aware SSL-based MOS predictor uses a pitch histogram as a conditioning vector for plain SSL-based MOS prediction framework. Since the pitch histogram statistically represents the pitch distribution of an utterance, referring to [18], it is concatenated with the utterance-level features obtained by pooling the frame-level features outputted by the SSL model. This is different from the aforementioned compressed-pitch-aware SSL-based MOS predictor. Then we balance these two types of features using a layer normalization and finally output the 1-dimensional MOS score through a linear layer.

3.2. Spectrum-aware SSL-based MOS predictor

As mentioned in Section 2.3, neural audio codecs differ from SSL models in that they can extract acoustic features from audio waveforms, and acoustic information is crucial for singing quality assessment. We use APCodec [8] to extract spectral-level acoustic features because it uses the audio amplitude and phase spectra as coding objects. In our preliminary experiments, we found that the quantization process affects the quality of the extracted features. While quantization is unavoidable for audio compression, it is not necessary for quality assessment. Therefore, we discard the quantizer and construct a non-quantized APCodec, consisting only of an encoder and a decoder. As depicted in Figure 3(c), for spectrum-aware SSL-based MOS predictor, a well-trained APCodec encoder extracts spectral-level acoustic features from singing voice waveforms. The frame shift used in this operation is the same as that used in the SSL model. Following this, a Conformer [19] is leveraged to capture a global view of the spectral representations from the input acoustic features. Similar to compressed-pitch-aware SSL-based MOS predictor, the output of the Conformer is concatenated with the output of the SSL model along the dimension axis, and then passed through a pooling layer and a linear layer to produce the 1-dimensional MOS score.

3.3. Bias Correction Branch

We propose a bias correction branch to replace the linear layer that outputs the MOS score in SSL-based MOS predictors, aiming at addressing the issue of imbalanced training sample distribution. As shown in Figure 4, we conducted a quantitative analysis of the

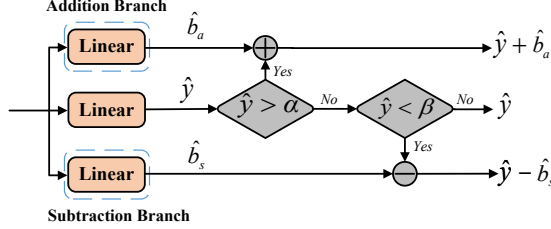


Fig. 5: The structure of the bias-correction branch.

SingMOS dataset provided by track 2 of VoiceMOS Challenge 2024 based on the score intervals of the labels. It can be observed that there are few samples in the low MOS intervals in both the training and validation sets, with the majority of training samples concentrated in the MOS interval above 3.0. Some other datasets [10] have MOS scores concentrated in the middle range, lacking high-MOS samples. We hypothesize that this imbalance in the training sample distribution may cause the trained model to have weaker predictive capabilities for samples with labels in certain MOS ranges compared to others, resulting in significant deviations between the predicted MOS and the actual labels.

To address this, we design a bias correction branch composed of three parallel linear layers that can theoretically be applied to any MOS prediction model by attaching it to the model’s output. As shown in Figure 5, the main improvement of the bias correction branch, compared to the commonly used linear layer at the output end of SSL-based MOS predictors, is the introduction of addition branch and subtraction branch. One branch is responsible for correcting the predicted scores of high MOS samples, while the other branch corrects the predicted scores of low MOS samples. Assume that the middle linear layer outputs the original MOS score \hat{y} , while addition branch and subtraction branch output bias values \hat{b}_a and \hat{b}_s , respectively. Given MOS score thresholds α and β ($1 < \beta < \alpha < 5$), the final output MOS score \hat{y}_{bc} of the bias correction branch is defined as

$$\hat{y}_{bc} = \begin{cases} \hat{y} + \hat{b}_a, & \hat{y} > \alpha \\ \hat{y}, & \beta < \hat{y} < \alpha \\ \hat{y} - \hat{b}_s, & \hat{y} < \beta \end{cases} . \quad (6)$$

This approach can prevent the MOS predictor from being confined to the middle MOS score range, allowing it to break through the interval boundaries and improve its assessment capabilities for both high and low MOS score samples. In the actual training process, after finishing training the original SSL-based MOS predictor, addition and subtraction branches are introduced at the model’s output, and only these two branches are trained while keeping the other model parameters fixed.

3.4. Model Fusion

According to previous studies [20, 21, 22, 23], it is evident that different types of MOS predictors can capture various information present in the training set. Integrating these models can effectively leverage these diverse perspectives to enhance the predictive capabilities of the quality assessment system. Therefore, we employ multiple pre-trained SSL models from *fairseq* [24] as our SSL models in plain, pitch-aware and spectrum-aware SSL-based MOS predictors. Then, we rank all MOS predictors according to preset rules and select the top predictors. After introducing the bias correction branches to each selected predictor, we perform the model fusion operation. Specifically, we concatenate the outputs of selected predictors into a

single fusion vector as input features, then a linear layer acts as the combiner to output the final fused MOS score.

4. EXPERIMENTS

4.1. Dataset and Evaluation Metrics

During the training phase, the track 2 of the VoiceMOS Challenge 2024 released a dataset named SingMOS. This dataset comprises Mandarin and Japanese samples obtained from SVS systems, SVC systems, analysis-synthesis operation of neural vocoders, and natural singing voice recordings. In total, the dataset includes 3189 samples, with 2000 samples allocated to the training set, 544 samples to the validation set, and 645 samples to the evaluation set. Notably, the validation and evaluation sets include unseen samples.

The evaluation metrics encompass mean squared error (MSE), linear correlation coefficient (LCC), SRCC, and kendall tau rank correlation (KTAU) at both the utterance and system levels [7]. MSE quantifies the disparity between predicted and ground truth MOS scores, whereas the other metrics assess correlation. According to the rules of the track 2 of VoiceMOS Challenge 2024, the system-level SRCC is used as the ranking criterion.

4.2. Comparison among Different SSL-based MOS Predictors

First, we compared the performance of the proposed pitch-aware and spectrum-aware SSL-based MOS predictors as well as the plain SSL-based MOS predictor to validate the effectiveness of the introduced pitch and spectral information for singing quality assessment. We selected 5 pre-trained models to use in these predictors, including Wav2Vec2.0 Base, Wav2Vec2.0 Large, HuBERT Base, HuBERT Large and HuBERT Extra Large. Therefore, a total of 20 predictors were compared, as shown in Table 1. For the spectrum-aware SSL-based MOS predictor, we adopted a no-quantized APCodec pre-trained on the VCTK-0.92 corpus [25] and a 2-layer Conformer, and the output feature dimension of spectral-level acoustic features was 64. During training, we trained all predictors for 1,000 epochs, with a batch size of 4 and optimizer as stochastic gradient descent (SGD) with a learning rate of 0.0001. For the checkpoint saving strategy, we followed the same approach as UTMOS [26], selecting the best system-level SRCC checkpoint calculated from the validation set. If the system-level SRCC didn’t decrease within 15 epochs, early stopping was applied. The purpose of comparing these predictors, besides validating the effectiveness of the introduced pitch and spectral information, is to select models for fusion. Hence, the experiments were conducted on the validation set.

The experimental results are shown in Table 1. For the two pitch-aware SSL-based MOS predictors, the pitch-histogram-aware SSL-based MOS predictor (i.e., PH-SSL-MOS in Table 1) demonstrated more stable performance. Compared to the plain SSL-based MOS predictor (i.e., SSL-MOS in Table 1), the pitch-histogram-aware one showed significant improvements in most metrics regardless of the SSL model used. Unfortunately, the compressed-pitch-aware SSL-based MOS predictor (i.e., CP-SSL-MOS in Table 1) seems to be sensitive to the SSL model. Specifically, when using Wav2Vec2.0 Large, its performance at both the utterance and system levels is very poor. The experimental results above confirm the effectiveness of introducing pitch-related information for singing quality assessment. Among the methods, the pitch histogram is more suitable for this task compared to pitch values. This may be attributed to the pitch histogram’s better representation of the characteristics of singing voices, confirming the hypothesis in Section 2.2. For the

Table 1: Experimental results of plain SSL-based MOS predictor (SSL-MOS), compressed-pitch-aware SSL-based MOS predictor (CP-SSL-MOS), pitch-histogram-aware SSL-based MOS predictor (PH-SSL-MOS) and spectrum-aware SSL-based MOS predictor (S-SSL-MOS) when adopting different SSL models evaluated on the validation set of SingMOS. The **bold** and underline numbers indicate the optimal and sub-optimal results, respectively.

Predictor	SSL Model	Utterance-level Metrics				System-level Metrics			
		MSE ↓	LCC ↑	SRCC ↑	KTAU ↑	MSE ↓	LCC ↑	SRCC ↑	KTAU ↑
SSL-MOS	Wav2Vec2.0 Base	0.470	0.652	0.647	0.487	<u>0.126</u>	0.945	0.928	0.775
CP-SSL-MOS		<u>0.370</u>	<u>0.668</u>	<u>0.635</u>	<u>0.475</u>	0.036	0.952	0.939	0.813
PH-SSL-MOS		0.344	0.669	0.629	0.474	0.241	<u>0.951</u>	0.939	<u>0.806</u>
S-SSL-MOS		0.477	0.660	0.623	0.465	0.159	0.908	<u>0.931</u>	0.783
SSL-MOS	Wav2Vec2.0 Large	<u>0.458</u>	<u>0.629</u>	<u>0.604</u>	<u>0.448</u>	<u>0.117</u>	0.939	<u>0.922</u>	<u>0.768</u>
CP-SSL-MOS		0.654	0.261	0.281	0.199	0.325	0.603	0.606	0.445
PH-SSL-MOS		0.320	0.689	0.651	0.489	0.035	<u>0.935</u>	0.933	0.787
S-SSL-MOS		0.543	0.501	0.469	0.334	0.202	0.780	0.749	0.571
SSL-MOS	HuBERT Base	0.406	0.642	0.619	0.459	0.077	0.925	0.910	0.749
CP-SSL-MOS		<u>0.383</u>	0.612	0.611	0.450	0.036	0.932	0.921	0.775
PH-SSL-MOS		0.338	0.685	0.653	0.489	<u>0.051</u>	0.949	0.935	0.806
S-SSL-MOS		0.966	<u>0.656</u>	<u>0.627</u>	<u>0.465</u>	0.620	<u>0.945</u>	<u>0.931</u>	<u>0.798</u>
SSL-MOS	HuBERT Large	0.485	0.533	0.524	0.378	0.093	0.823	0.855	0.688
CP-SSL-MOS		0.434	0.553	0.539	0.397	<u>0.092</u>	0.874	0.876	0.715
PH-SSL-MOS		0.399	0.588	<u>0.573</u>	0.420	0.059	0.904	0.894	<u>0.726</u>
S-SSL-MOS		1.387	0.616	0.586	0.429	1.026	<u>0.886</u>	0.896	0.730
SSL-MOS	HuBERT Extra Large	0.443	0.543	0.539	0.395	0.064	0.876	0.887	0.715
CP-SSL-MOS		<u>0.421</u>	0.602	0.599	0.440	<u>0.048</u>	0.909	<u>0.918</u>	0.765
PH-SSL-MOS		0.359	0.651	0.641	0.478	0.043	<u>0.919</u>	0.923	<u>0.772</u>
S-SSL-MOS		1.873	<u>0.635</u>	<u>0.600</u>	<u>0.442</u>	1.274	0.936	0.923	0.779

spectrum-aware SSL-based MOS predictor (i.e., S-SSL-MOS in Table 1), when using Wav2Vec2.0 Large, its performance is also poor, and regardless of the SSL model used, its MSE metric significantly increases compared to the plain predictor. Nevertheless, its performance is satisfactory regarding ranking metrics, such as system-level SRCC. Therefore, the introduction of spectral-level acoustic information is helpful in improving the accuracy of singing quality assessment models to a certain extent.

4.3. Validation on Model Fusion and Bias Correction

As mentioned in Section 3.4, the PS-SQA employed a model fusion strategy for better singing MOS prediction accuracy. We selected five predictors from the twenty predictors shown in Table 1 for fusion and constructed the PS-SQA. Since system-level SRCC is an important ranking metric in VoiceMOS Challenge 2024, we selected the top five predictors based on their system-level SRCC rankings. The selected predictors are listed in Table 2, including 1) compressed-pitch-aware MOS predictor with Wav2Vec2.0 Base as the SSL model, 2) pitch-histogram-aware MOS predictor with Wav2Vec2.0 Base as the SSL model, 3) pitch-histogram-aware MOS predictor with Wav2Vec2.0 Large as the SSL model, 4) pitch-histogram-aware MOS predictor with HuBERT Base as the SSL model, and 5) spectrum-aware MOS predictor with HuBERT Base as the SSL model. We can observe that among the five selected predictors, four incorporate pitch information, while only one incorporates spectrum information. This further indicates that pitch information is more important for enhancing the performance of singing quality assessment methods compared to spectrum information. This may be due to the specific characteristics of singing voices. Within the pitch-aware methods, three predictors based on pitch histograms were selected, while only one predictor based on compressed pitch was chosen. This indicates that pitch histograms

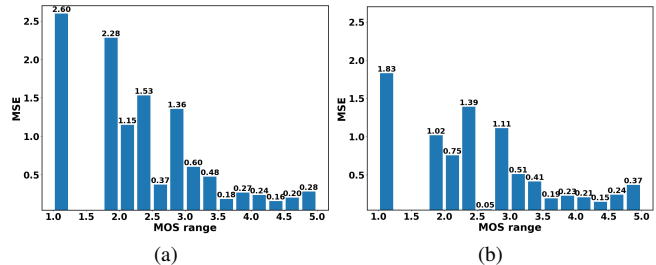


Fig. 6: Bar charts showing the MSE for each MOS score segment (interval length of 0.25) for pitch-histogram-aware MOS predictor with Wav2Vec2.0 Base as the SSL model, where subplot (a) displays the results without bias correction and subplot (b) displays the results with bias correction.

can provide more easily interpretable pitch information for SSL-based MOS prediction models, facilitating the learning process for the models. As the results of PS-SQA without bias correction in Table 2 (i.e., the rows annotated with X in the “Bias Correction” column) suggested, except for the system-level LCC metric, the performance of the fused PS-SQA is higher than any of the individual predictors in all other metrics.

Specifically, after model fusion, the system-level SRCC improved by 0.013 compared to the best individual predictor. This indicates that the model fusion strategy combines the strengths of individual MOS predictors, achieving more accurate singing MOS predictions. This confirms the effectiveness of model fusion.

Then, we added the bias correction branch to each of the individual predictors to explore the role of bias correction. The experimental results are also listed in Table 2. By comparing the performance of individual predictors with (i.e., the rows annotated with ✓ in the

Table 2: Experimental results of the fused model of PS-SQA and included predictors without or with bias correction branch evaluated on the validation set of SingMOS. The **bold** number indicates the optimal results.

Fusion Model & Predictor	SSL Model	Bias Correction	Utterance-level Metrics				System-level Metrics			
			MSE ↓	LCC ↑	SRCC ↑	KTAU ↑	MSE ↓	LCC ↑	SRCC ↑	KTAU ↑
PS-SQA		×	0.302	0.710	0.674	0.507	0.031	0.946	0.952	0.840
CP-SSL-MOS	Wav2Vec2.0 Base	×	0.370	0.668	0.635	0.475	0.036	0.952	0.939	0.813
PH-SSL-MOS	Wav2Vec2.0 Base	×	0.344	0.669	0.629	0.474	0.241	0.951	0.939	0.806
PH-SSL-MOS	Wav2Vec2.0 Large	×	0.320	0.689	0.651	0.489	0.035	0.935	0.933	0.787
PH-SSL-MOS	HuBERT Base	×	0.338	0.685	0.653	0.489	0.051	0.949	0.935	0.806
S-SSL-MOS	HuBERT Base	×	0.966	0.656	0.627	0.465	0.620	0.945	0.931	0.798
PS-SQA		✓	0.320	0.685	0.672	0.506	0.027	0.944	0.953	0.840
CP-SSL-MOS	Wav2Vec2.0 Base	✓	0.492	0.679	0.641	0.480	0.143	0.954	0.944	0.828
PH-SSL-MOS	Wav2Vec2.0 Base	✓	0.374	0.654	0.627	0.471	0.040	0.959	0.945	0.828
PH-SSL-MOS	Wav2Vec2.0 Large	✓	0.328	0.686	0.651	0.489	0.042	0.941	0.938	0.802
PH-SSL-MOS	HuBERT Base	✓	0.341	0.687	0.654	0.489	0.054	0.950	0.938	0.813
S-SSL-MOS	HuBERT Base	✓	0.477	0.660	0.623	0.465	0.159	0.908	0.931	0.793

Table 3: Experimental results of the official baseline, the participating teams, and our proposed PS-SQA on the test set of SingMOS. The **bold** and underline numbers indicate the optimal and sub-optimal results, respectively.

Participating Systems	Utterance-level Metrics				System-level Metrics			
	MSE ↓	LCC ↑	SRCC ↑	KTAU ↑	MSE ↓	LCC ↑	SRCC ↑	KTAU ↑
B01 (Official Baseline)	0.419	0.594	0.605	0.442	0.079	0.851	<u>0.859</u>	<u>0.687</u>
T01	0.366	0.605	0.603	0.440	<u>0.051</u>	0.858	0.837	0.684
T03	0.432	0.597	0.583	0.426	0.061	0.848	0.819	0.637
T04	<u>0.363</u>	0.624	0.604	0.445	0.056	<u>0.869</u>	0.833	0.640
T05	<u>0.363</u>	0.609	0.593	0.434	0.069	0.791	0.807	0.657
T06	0.358	0.637	<u>0.625</u>	<u>0.460</u>	0.063	0.841	0.831	0.657
T08 (Our Submission)	0.384	0.628	0.620	0.455	0.072	0.845	0.856	<u>0.687</u>
PS-SQA	0.452	<u>0.634</u>	0.639	0.470	0.031	0.880	0.888	0.717

“Bias Correction” column) and without the bias correction branch, it can be observed that the introduction of the bias correction branch effectively improved performance in most metrics. Figure 6 also shows the MSE metrics for pitch-histogram-aware MOS predictor with Wav2Vec2.0 Base as the SSL model across different MOS segments using bar charts, comparing the results without and with bias correction. The label MOS scores in the validation set, ranging from 1 to 5, are evenly divided into 16 segments with an interval of 0.25. There are no samples in the 2nd and 3rd MOS segments. By comparing the two subplots in Figure 4, it is clear that in the low MOS segments, the introduction of bias correction significantly reduced the MSE metric. In the MOS segment from 1.75 to 2, the MSE reduction even exceeded 1. This indicates that bias correction effectively enhances the predictor’s modeling capability on low-resource data, effectively mitigating the difficulty the predictor faces when handling specific MOS segments. When these five bias-corrected predictors are fused to construct the PS-SQA, the performance also surpassed that of the individual predictors. The bias-corrected fused PS-SQA also shows a slight improvement in system-level metrics compared to the non-bias-corrected fused PS-SQA.

4.4. Position in VoiceMOS Challenge 2024 Competition Systems

Based on the results discussed in Sections 4.2 and 4.3 on the validation set, we chose the PS-SQA with bias correction and model fusion strategies as our competitive system to compare with the systems participating in track 2 of VoiceMOS Challenge 2024. This allows us to evaluate the position of our proposed method among these competition systems. Table 3 presents the experimental results of the official baseline, the participating teams, and the PS-SQA proposed in this paper on the SingMOS evaluation set. Our submitted

system (i.e., T08) significantly outperformed the other participating systems (i.e., T01, T03, T04, T05 and T06) in terms of the system-level SRCC metric used for ranking, although it was slightly inferior to the official baseline. Our submitted system only utilized the pitch-histogram-aware SSL-based MOS predictor as shown in Figure 3(b) without layer normalization. Building on this, we made improvements and proposed other-information-aware SSL-based MOS predictors. We also introduced the bias correction branch in various predictors and adopted a model fusion strategy. Currently, in terms of system-level metrics, the proposed PS-SQA significantly outperformed all other systems. Specifically, compared to our submitted system, PS-SQA improved the system-level SRCC by 0.032.

5. CONCLUSION

This paper proposes a novel pitch-and-spectrum-aware singing quality assessment method, called PS-SQA, which is an improvement version of the system we submitted to track 2 of VoiceMOS Challenge 2024. The PS-SQA first introduces multiple MOS predictors that incorporate pitch and spectrum-related information into the SSL-based MOS prediction model, tailored to the characteristics of singing voices. The PS-SQA then selects these predictors for model fusion and introduces a bias correction branch to overcome training bias caused by low-resource training data. Experimental results confirm that our proposed PS-SQA significantly outperforms all participating systems in terms of system-level evaluation metrics. Enhancing the MOS prediction accuracy of PS-SQA by introducing more types of predictors will be our future work.

6. REFERENCES

- [1] Lin Huang, Chitralkha Gupta, and Haizhou Li, “Spectral features and pitch histogram for automatic singing quality evaluation with CRNN,” in *Proc. APSIPA*, 2020, pp. 492–499.
- [2] Jinhu Li, Chitralkha Gupta, and Haizhou Li, “Training explainable singing quality assessment network with augmented data,” in *Proc. APSIPA*, 2021, pp. 904–911.
- [3] Xiaoheng Sun, Yuejie Gao, Hanyao Lin, and Huaping Liu, “Tg-Critic: A timbre-guided model for reference-independent singing evaluation,” in *Proc. ICASSP*, 2023, pp. 1–5.
- [4] A Black and Keiichi Tokuda, “The Blizzard Challenge 2005: Evaluating corpus-based speech synthesis on common databases,” in *Proc. interspeech*, 2005, pp. 77–80.
- [5] Brian Patton, Yannis Agiomyrgiannakis, Michael Terry, Kevin Wilson, Rif A Saurous, and D Sculley, “AutoMOS: Learning a non-intrusive assessor of naturalness-of-speech,” *arXiv preprint arXiv:1611.09207*, 2016.
- [6] Szu-wei Fu, Tsao Yu, Hsin-Te Hwang, and Hsin-Min Wang, “Quality-Net: An end-to-end non-intrusive speech quality assessment model based on BLSTM,” in *Proc. Interspeech*, 2018, pp. 1873–1877.
- [7] Chen-Chou Lo, Szu-Wei Fu, Wen-Chin Huang, Xin Wang, Junichi Yamagishi, Yu Tsao, and Hsin-Min Wang, “MOSNet: Deep learning-based objective assessment for voice conversion,” in *Proc. Interspeech*, 2019, pp. 1541–1545.
- [8] Yang Ai, Xiao-Hang Jiang, Ye-Xin Lu, Hui-Peng Du, and Zhen-Hua Ling, “APCodec: A neural audio codec with parallel amplitude and phase spectrum encoding and decoding,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 32, pp. 3256–3269, 2024.
- [9] Erica Cooper, Wen-Chin Huang, Tomoki Toda, and Junichi Yamagishi, “Generalization ability of MOS prediction networks,” in *Proc. ICASSP*, 2022, pp. 8442–8446.
- [10] Wen-Chin Huang, Erica Cooper, Yu Tsao, Hsin-Min Wang, Tomoki Toda, and Junichi Yamagishi, “The VoiceMOS Challenge 2022,” in *Proc. Interspeech*, 2022, pp. 4536–4540.
- [11] Erica Cooper, Wen-Chin Huang, Yu Tsao, Hsin-Min Wang, Tomoki Toda, and Junichi Yamagishi, “The VoiceMOS Challenge 2023: zero-shot subjective speech quality prediction for multiple domains,” in *Proc. ASRU*, 2023, pp. 1–7.
- [12] George Tzanetakis, Andrey Ermolinskyi, and Perry Cook, “Pitch histograms in audio and symbolic music information retrieval,” *Journal of New Music Research*, pp. 143–152, 2003.
- [13] Chitralkha Gupta, Haizhou Li, and Ye Wang, “Perceptual evaluation of singing quality,” in *Proc. APSIPA*. IEEE, 2017, pp. 577–586.
- [14] Chitralkha Gupta, Haizhou Li, and Ye Wang, “Automatic evaluation of singing quality without a reference,” in *Proc. APSIPA*, 2018, pp. 990–997.
- [15] Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi, “Soundstream: An end-to-end neural audio codec,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 495–507, 2021.
- [16] Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi, “High fidelity neural audio compression,” *Transactions on Machine Learning Research*, 2023.
- [17] Ryandhimas E Zezario, Szu-Wei Fu, Fei Chen, Chiou-Shann Fuh, Hsin-Min Wang, and Yu Tsao, “Deep learning-based non-intrusive multi-objective speech assessment model with cross-domain features,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 31, pp. 54–70, 2022.
- [18] Chitralkha Gupta, Lin Huang, and Haizhou Li, “Automatic rank-ordering of singing vocals with twin-neural network,” in *Proc. ISMIR*, 2020, pp. 416–423.
- [19] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al., “Conformer: Convolution-augmented transformer for speech recognition,” in *Proc. Interspeech*, 2020, pp. 5036–5040.
- [20] Marie Kunešová, Jindřich Matoušek, Jan Lehečka, Jan Švec, Josef Michálek, Daniel Tihelka, Martin Bulín, Zdeněk Hanzlíček, and Markéta Řezáčková, “Ensemble of deep neural network models for MOS prediction,” in *Proc. ICASSP*, 2023, pp. 1–5.
- [21] Zhengdong Yang, Wangjin Zhou, Chenhui Chu, Sheng Li, Raj Dabre, Raphael Rubino, and Yi Zhao, “Fusion of self-supervised learned models for MOS prediction,” in *Proc. Interspeech*, 2022, pp. 5443–5447.
- [22] Adriana Stan, “The ZovoMOS entry to VoiceMOS Challenge 2022,” in *Proc. Interspeech 2022*, 2022, pp. 4516–4520.
- [23] Zili Qi, Xinhui Hu, Wangjin Zhou, Sheng Li, Hao Wu, Jian Lu, and Xinkang Xu, “LE-SSL-MOS: Self-supervised learning MOS prediction with listener enhancement,” in *Proc. ASRU*, 2023, pp. 1–6.
- [24] Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli, “Fairseq: A fast, extensible toolkit for sequence modeling,” *arXiv preprint arXiv:1904.01038*, 2019.
- [25] Junichi Yamagishi, Christophe Veaux, Kirsten MacDonald, et al., “CSTR VCTK corpus: English multi-speaker corpus for CSTR voice cloning toolkit (version 0.92),” *University of Edinburgh. The Centre for Speech Technology Research (CSTR)*, 2019.
- [26] Takaaki Saeki, Detai Xin, Wataru Nakata, Tomoki Koriyama, Shinnosuke Takamichi, and Hiroshi Saruwatari, “UTMOS: UTokyo-SaruLab system for VoiceMOS Challenge 2022,” in *Proc. Interspeech*, 2022, pp. 4521–4525.