

Bias-Aware Loss for Training Image and Speech Quality Prediction Models from Multiple Datasets

Gabriel Mittag¹, Saman Zadtootaghaj¹, Thilo Michael¹, Babak Naderi¹, Sebastian Möller^{1,2}

¹Quality and Usability Lab, Technische Universität Berlin, Berlin, Germany

²Deutsches Forschungszentrum für Künstliche Intelligenz (DFKI), Berlin, Germany
first.last@tu-berlin.de

Abstract—The ground truth used for training image, video, or speech quality prediction models is based on the Mean Opinion Scores (MOS) obtained from subjective experiments. Usually, it is necessary to conduct multiple experiments, mostly with different test participants, to obtain enough data to train quality models based on machine learning. Each of these experiments is subject to an experiment-specific bias, where the rating of the same file may be substantially different in two experiments (e.g. depending on the overall quality distribution). These different ratings for the same distortion levels confuse neural networks during training and lead to lower performance. To overcome this problem, we propose a bias-aware loss function that estimates each dataset’s biases during training with a linear function and considers it while optimising the network weights. We prove the efficiency of the proposed method by training and validating quality prediction models on synthetic and subjective image and speech quality datasets.

Index Terms—Speech Quality, Image Quality, DNN

I. INTRODUCTION

In order to optimise the Quality of Experience (QoE) of multimedia services, developers rely on measures to validate new codecs or communication channels. Traditionally, the quality of image, video, or speech services is measured in subjective experiments, in which test participants are asked to rate the quality of a sample. The average across all test participants’ ratings gives the so-called Mean Opinion Score (MOS). However, because this procedure is time-consuming and costly, instrumental models have been developed to automatically estimate the quality. In the case of speech quality, full-reference models, such as PESQ and POLQA have been established. In the image quality domain, many different visual quality metrics have been proposed [1]. The most popular metrics are NIQE [2] and BRISQUE [3] for no-reference assessment or PSNR and SSIM [4] for full-reference assessment. More recently also deep learning approaches have been introduced for speech quality or synthesised speech naturalness [5]–[9] and for image quality [10]–[12] prediction.

While objective quality models always compute the same score for a given sample, the MOS determined through subjective quality experiments, which is used as ground truth for training such objective models, is a sensitive measure. Minor change in a vote given from one test participant leads to a change of the overall value [13]. Consequently, test-retest studies even with the same group of participants, who rate the same dataset, often do not lead to the exact same MOS

values [13], [14]. Given that, it is recommended to consider each subjective experiment as a closed set [14].

Datasets with subjective MOS are usually limited in size due to the maximum time in which one participant is able to rate the corpus before fatigue occurs. Therefore, it is common practice to use multiple datasets for training deep-learning-based quality prediction models. Furthermore, subjective data is usually sparse due to the costs that experiments involve and thus also older datasets are included for model training to increase the training size. Consequently, as these datasets often come from different labs with different test participants, and often, many years lie in between them, they are exposed to dataset-specific, subjective biases that make a comparison of the MOS values from different experiments difficult. The main bias-inducing factors according to ITU-T Rec. P.1401 [14] are:¹

- **Rating noise** The score assigned by a listener is not always the same, even if an experiment is repeated with the same samples and the same presentation order.
- **Order-effect** Subjects are influenced by the short-term history of the samples they previously rated. For example, after one or two poor samples, participants tend to rate a mediocre sample higher. In contrast, if a mediocre sample follows high-quality samples, there is a tendency to score the mediocre sample lower. Because of this effect, the presentation order for each subject is usually randomised in quality experiments.
- **Corpus-effect** The largest influence is given by effects associated with the average quality, the distribution of quality, and the occurrence of individual distortions. Test participants tend to use the entire set of scores offered in an experiment. Because of this, in an experiment that contains mainly low-quality samples, the subjects will overall rate them higher, introducing a constant bias. Despite verbal category labels, subjects adapt the scale to the qualities presented in each experiment. Furthermore, individual distortions that are presented less often are rated lower, as compared to experiments in which samples are presented more often, and people become more familiar with them. For example, it was shown in [16] that a clean narrowband speech signal obtains a higher MOS in a corpus with only narrowband conditions than in

¹Besides these factors, [15] lists further biases that occur in listening tests.

a mixed-band corpus that includes wideband conditions as well.

- **Long-term dependencies** These effects reflect the general cultural behaviour of the subjects as to the exact interpretation of the category labels, the cultural attitude to quality, and language dependencies. Also, the daily experiences with telecommunication or media are important. Quality experience, and therefore expectation, may change over time as the subjects become more familiar with high-definition codecs and VoIP distortions.

In this paper, *bias* refers to offsets and gradients between datasets that are caused by these factors. While offsets are mostly introduced by the overall quality that is presented to the participants during an experiment, gradients are, for example, introduced when the experiment does not cover the entire quality range (i.e. the ratings tend to become more pessimistic faster).

These induced biases result in different MOS values for stimuli of the same distortion levels if they are contained in a different dataset. As a consequence, when multiple datasets are combined for training a quality prediction model, the correct rank order of MOS values is no longer given. The error-prone rank order of the combined training set leads to lower prediction performance of the trained quality prediction model. To overcome this problem, usually a set of common anchor conditions are included in all datasets (typically 20% of test conditions). These anchor conditions cover the whole range of quality distortions and make a comparison of different datasets possible. By calculating a mapping function between the anchor conditions of the different datasets, biases can be partly removed before training. However, if the datasets that are used for training originate from different sources, they often do not contain the same anchor conditions. It should further be noted that biases that are introduced by the corpus-effect cannot be removed by normalising the datasets, which is a common preprocessing step. In case two datasets contain a different range of distortion levels, normalisation of the MOS values may even increase the bias between datasets.

So far in literature, all deep learning based image or speech quality prediction models (e.g. the aforementioned models in [5]–[12]) do not consider these biases between datasets and apply either a vanilla MSE (mean squared error) or MAE (mean absolute error) loss function.

In this paper, we present a method that deals with subjective biases between datasets in the training phase of deep neural network models without the need for anchor conditions. The novelty of the proposed loss function is that it learns the biases automatically by using the model predictions as objective, unbiased measure. The proposed bias-aware loss is repeatedly updated during the model training, which avoids an overfitting to the biases themselves. The presented algorithm in this paper together with the synthetic bias experiments are made publicly available.²

The rest of the paper is structured as follows: At first we present the proposed loss function and the corresponding algorithm. The method is then evaluated on a synthetic dataset for which the biases are known. After that the bias-aware loss is applied to speech and image quality datasets with subjective MOS ratings.

II. METHOD

The basic idea is that the dataset-specific biases are learned and accounted for during the training of the neural network. After each epoch, the predicted MOS values of the training data are used to estimate the bias in each dataset by mapping their values to the ground truth subjective MOS values. The hypothesis is that the predicted MOS values of the model are objective in the sense that they will average out the biases of the different datasets while training.

The biases are approximated with a first-order polynomial function. After the biases are estimated for each dataset, they can be used in the next epoch to calculate the proposed bias-aware loss. To this end, the predicted MOS values are mapped with the calculated bias coefficients. The MSE between the bias-mapped predicted MOS and the subjective MOS values then gives the loss as follows:³

$$l = \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}, \mathbf{b}) = \frac{1}{N} \sum_{i=1}^N \left(y_i - (b_0^j + b_1^j \hat{y}_i) \right)^2 \quad (1)$$

where i is the index of the sample belonging to dataset j , y_i is the subjective MOS and \hat{y}_i is the predicted value for that sample, b_0^j and b_1^j are the estimated bias coefficients for dataset j , and N is the overall number of training samples.

Because the predicted values are mapped according to the bias of each dataset, errors between the predicted and subjective values that only occur due to the dataset-specific bias are neglected. The model thus learns to predict quality, rather than non-relevant biases.

A. Learning with bias-aware loss

The complete algorithm of the bias-aware loss is depicted in Algorithm 1. The inputs to the algorithm are the input features x_i and the subjective MOS values that are the desired output values y_i for all samples $i \in N$. Further, a list that contains the dataset index j of each sample db_i is needed to assign the individual samples to their datasets. Before the training starts, the bias coefficients b^j of all datasets are initialised with an identity function, with which Eq. (1) corresponds to a vanilla MSE loss. Because the model will typically not give a meaningful prediction output after the first few epochs, the `update_bias` flag is set to `False` until a predefined model accuracy r_{th} in terms of Pearson’s correlation coefficient (PCC) is achieved (see analysis in Sect III-A). Until this threshold is not reached the bias coefficients will not be updated, and therefore, a vanilla MSE loss is used for calculating the loss. As a metric PCC is used instead of RMSE (root-mean-square

²<https://github.com/gabrielmittag/Bias-Aware-Loss>

³In this paper, we use MSE as the base loss function, however, the bias-aware loss algorithm can be applied to any distance measurement function.

error) as the RMSE is strongly affected by biases and therefore unsuitable in this case.

Algorithm 1 Training with bias-aware loss function

Input: x_i : input features, y_i : subjective MOS values, db_i : list of dataset indices j for each sample i

Parameter: r_{th} : minimum prediction accuracy to update the bias coefficients

Output: model weights

```

1: Number of dataset:  $D = \max(\mathbf{db})$ 
2: Initialise bias for each dataset:  $b^j = [0, 1], j \in D$ 
3: update_bias = False
4: while not converged do
5:   Shuffle mini-batch indices  $idx_k$ 
6:    $k = 0$ 
7:   for all mini-batches do
8:     Get mini-batch:
        $\mathbf{x}_b = \mathbf{x}[idx_k], \mathbf{y}_b = \mathbf{y}[idx_k]$ 
9:     Feed forward:  $\hat{\mathbf{y}}_b = \text{model}(\mathbf{x}_b)$ 
10:    Calculate bias-aware loss with Eq. (1):
       $l = \mathcal{L}(\mathbf{y}_b, \hat{\mathbf{y}}_b, b^j)$ .
11:    Backpropagate & optimise weights
12:     $k = k + 1$ 
13:  end for
14:  Predict MOS:  $\hat{\mathbf{y}} = \text{model}(\mathbf{x})$ 
15:  Calculate Pearson's correlation  $r = \text{PCC}(\mathbf{y}, \hat{\mathbf{y}})$ .
16:  if  $r > r_{th}$  or update_bias then
17:    update_bias = True
18:    for  $j$  in  $D$  do
19:      Find dataset indices:
         $idx = \text{find}(\mathbf{db} == j)$ 
20:      Get dataset:  $\mathbf{y}^{db} = \mathbf{y}[idx], \hat{\mathbf{y}}^{db} = \hat{\mathbf{y}}[idx],$ 
         $\mathbf{M} = \text{len}(idx)$ 
21:      Estimate bias:
         $\min_{b^j} \frac{1}{M} \sum_{i=1}^M (y_i^{db} - (b_0^j + b_1^j \hat{y}_i^{db}))$ 
22:    end for
23:  end if
24: end while
25: return model weights

```

At the start of each epoch, the mini-batch indices idx_k are randomly shuffled. It is necessary to preserve these indices in order to assign the samples to their corresponding datasets. After each epoch, the model is used to predict the MOS values $\hat{\mathbf{y}}$ of all samples. Once the model accuracy r_{th} is reached, the update_bias flag is set to True and the biases will be estimated after every epoch.

To update the bias coefficients, the algorithm loops through all datasets individually, where idx represents the indices of all samples that belong to the dataset j . At line 20 of the algorithm, the subjective and predicted MOS values of the samples belonging to dataset j are loaded. Then they are used to estimate the bias coefficients b^j . In the next training epoch, the biases are then applied to calculate the loss (line 10) of each mini-batch that may contain a random number

of different datasets, and therefore each sample may also be subject to a different bias. By using the bias-aware loss, these biases are considered when calculating the error between the predicted and subjective MOS values. After each epoch, the bias coefficients are then updated to be in line with the updated model predictions. For more details on the algorithm see also the open-sourced PyTorch code.

B. Anchoring Predictions

When the bias-aware loss is applied, the MOS predictions are not anchored to the absolute subjective MOS values and, therefore, can wander off. While the predictions will still rank the samples in the best possible way, there may be a large offset between predictions and subjective values on the validation set. This effect will lead to a higher RMSE while the PCC is usually not affected. To overcome this problem, the predicted MOS values of all samples can be mapped to the subjective MOS after the training. This mapping can then be applied when making new predictions on validation data.

Alternatively, instead of estimating the bias for all datasets, the predictions can be anchored to one specific training dataset. This approach can be particularly useful if there is one dataset of which it is known that the conditions are similar to the conditions that the model should be applied to later. A new dataset is usually created with new conditions and then split into a training and validation set. To increase the training data size and to improve the model accuracy, older datasets may also be included in the training set. It is likely that there will be a bias between the new dataset and the older dataset, for example, because the highest quality in the new dataset will be higher than the one in the older datasets. The predictions can be anchored to the new dataset by omitting the bias update for the new dataset only (skipping line 19-21 in Algorithm 1). The biases of all other datasets are then computed in relation to the anchor dataset and, as a result, the final model predictions will be in line with the anchor dataset.

III. EXPERIMENTS AND RESULTS

In the first experiment, we generate a synthetic speech quality dataset for which the biases are known. After that, we conduct two more experiments with real speech quality and image quality data. Because the results of each neural network training run depend on random initialisation, random shuffling and other factors, such as random dropout, we run each experiment 15 times and use the average results to rule out any random effects.

A. Synthetic data

Firstly, a synthetic speech quality dataset is generated to which artificial biases are applied. As source files, the 2–3 s reference speech files from the TSP dataset [17] are used. For the four training datasets overall 320 speech files and for the validation set 80 speech files are used. The speech signals were processed with white Gaussian noise to create conditions of different distortion levels. It was found that when the SNR range of the added noise is too wide, it is too easy for the

model to predict the quality, and when it is too narrow, the prediction becomes too challenging. Therefore, the speech files were processed with noise at SNR values between 20 dB and 25 dB, which showed to be a good compromise. To simulate MOS predictions, an S-shaped mapping between the technical impairment factor (i.e., SNR) and MOS, taken from ITU-T Rec. G.107 [18] is used. The relationship is shown in Figure 1 and maps the SNR values to a MOS between 1 and 4.5.

The training data is divided into four different subsets, and a different bias is applied to each of them. These four artificial biases are shown in Figure 2. To better analyse the influence of the bias-aware loss, extreme biases are used in this experiment. There is no bias applied to the first simulated database (blue line). The second and the third simulated databases have linear biases applied, while the fourth database is exposed to a bias modelled with a third-order polynomial function. Each of the training datasets contains 80 files; the validation set also contains 80 files. The synthetic experiments are run by using the speech quality model NISQA [19].

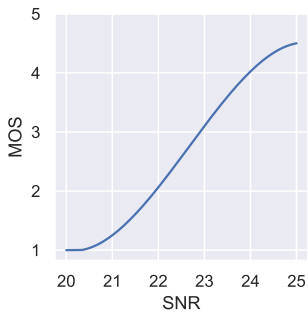


Fig. 1: Mapping between noise in SNR and speech quality MOS.

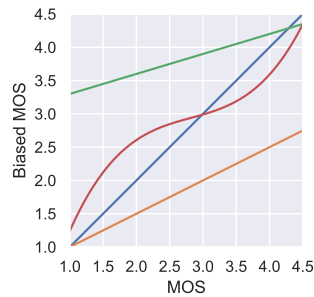


Fig. 2: Four artificial biases introduced to the four simulated training datasets.

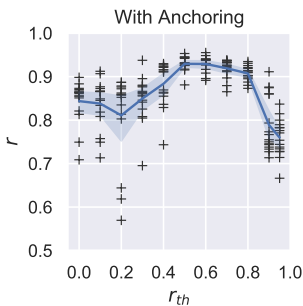


Fig. 3: Validation results in terms of Pearson's correlation over 15 training runs with anchoring for different r_{th} thresholds.

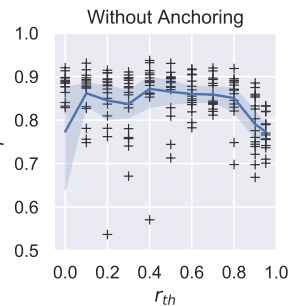


Fig. 4: Validation results in terms of Pearson's correlation over 15 training runs without anchoring for different r_{th} thresholds.

Minimum Accuracy r_{th} : In the first experiment, the influence of the minimum PCC that must be achieved on the training set before activating the bias update in the bias-aware loss algorithm (line 16 in Algorithm 1) is analysed. To this end,

the experiment is run with 11 different threshold r_{th} between 0 to 0.95. An early stop on the validation PCC of 20 epochs is used, and the best epoch of each run is saved as result. The training run of each of these 11 configurations is repeated 15 times. The results, together with the mean results and their 95% confidence interval, can be seen in Figure 4 for training without anchoring and in Figure 3 with anchoring. In the case of anchoring, the bias-aware loss is not used on the first dataset *train_1* but only on the other three datasets.

The correlation of the results is highly varying when an anchor dataset is used. The highest correlation can be achieved for thresholds between 0.5 and 0.7. When no anchoring is applied, the exact threshold does not seem to be as crucial, as long as it is somewhere between 0.1 and 0.8. However, the PCC remains overall lower than the higher PCCs that can be achieved with an anchor dataset. For a threshold higher than 0.9, the accuracy notably drops. It can be assumed that at this point, the model weights are already optimised too far towards the vanilla MSE loss and cannot always profit from the late activated bias-aware loss.

Synthetic training example: Figure 5 shows an example of four different training runs and different anchoring configuration. The figures show the epoch with the best results on the validation data set in terms of PCC. Each row presents one training run, and each column presents the results on the four different training datasets and on the validation dataset. The artificial bias that was applied to the datasets can be seen as green line (see also Figure 2). The estimated bias used by the bias-aware loss is depicted as orange line. The top row presents the results without anchoring and shows how the prediction results can drift away from the original values. While the predictions are extremely biased in this case, the achieved PCC remains high.

The second row shows the results for anchoring with anchoring dataset. During the training, the bias of the first training dataset *train_1* was not estimated but fixed to an identity function. It can be seen that the prediction results on the validation set are less biased in this case. Furthermore, it can be seen that the model successfully learns the different biases when the estimated orange line is compared to the original green line.⁴

Results on synthetic data The results on the validation dataset of 15 different runs are presented as boxplots in Figure 6. When the second boxplot is compared to the first one where the original, unbiased data was used for training, it can be seen that the model accuracy decreases significantly from an average PCC of $r = 0.95$ to $r = 0.77$ when the training datasets are exposed to biases. This performance decrease caused by the biased training data can successfully be compensated for by the bias-aware loss, which is displayed in the third boxplot on the right-hand side. The average PCC obtained is $r = 0.93$, which is close to the results trained on

⁴We also experimented with third-order polynomial estimation functions with which biases such as in dataset *train_4* can be estimated more precisely. However, because they did not further improve the results they are left out of this work.

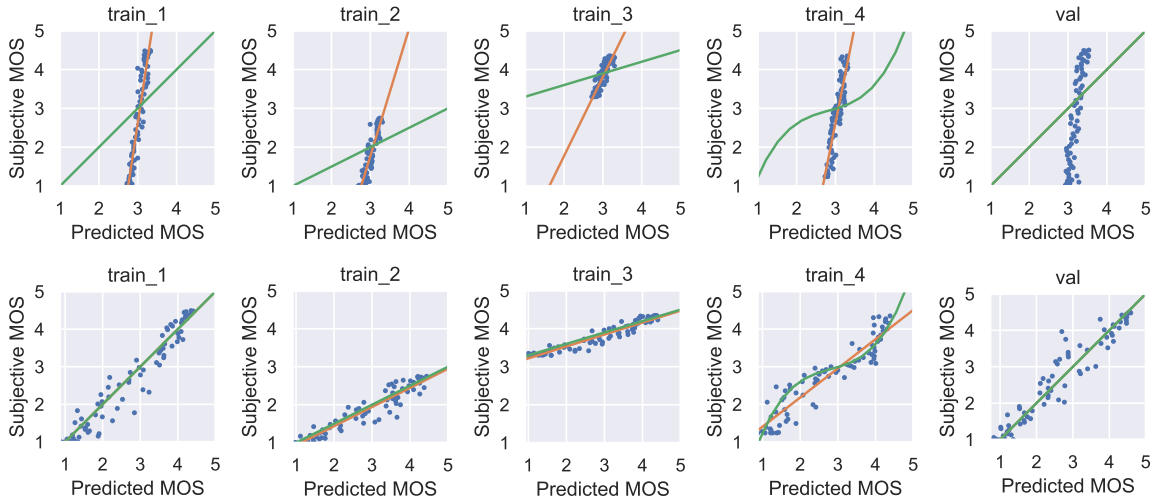


Fig. 5: Two example training runs with bias-aware loss on the validation dataset. GREEN LINE: Introduced bias. ORANGE LINE: Estimated bias. 1st Row: Without anchoring. 2nd Row: Anchored with dataset train_1.

unbiased data. Overall, the experiments show the efficiency of the proposed bias-aware loss when applied to the synthesised data, where the PCC can be increase from 0.77 to 0.93.

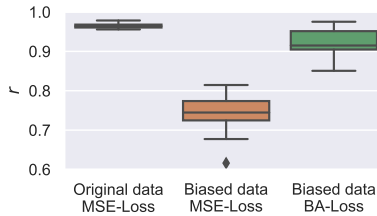


Fig. 6: Validation results of 15 training runs on the synthesised data.

B. Speech Quality

To analyse the efficiency on real datasets, the publicly available ITU-T Rec. P.Sup23 [20] speech quality datasets with subjective quality ratings are used. The datasets originate from different sources and are therefore likely to be exposed to subjective biases. Six of the datasets that were rated on an ACR scale are used: EXP1a, EXP1d, EXP1o, EXP3a, EXP3c, and EXP3d. The proposed algorithm is analysed with a leave-one-dataset-out cross-validation, where the model is trained on the five remaining datasets. For each training run, the model of the epoch on which the best results in terms of PCC was achieved on the held-out validation dataset is saved. As prediction model again, the speech quality model NISQA is applied.

The training was run 15 times with an early stop of 20 epochs on the validation PCC, a learning rate of 0.001, and a mini-batch size of 60. The bias estimation was anchored to a randomly chosen dataset. The hyper-parameter r_{th} was optimised for each individual dataset. The average results for each dataset over all runs are shown in Table I. It can be noted

that the efficiency of the proposed algorithm depends on the datasets it has been trained and evaluated on. For five of the six datasets, the proposed loss outperforms the vanilla MSE loss function, whereas, for dataset EXP1o, the vanilla MSE loss performs better. On the EXP3a dataset, a performance increase of about 0.05 in terms of PCC can be observed, showing that the speech quality prediction can notably be improved by applying the proposed bias-aware loss.

C. Image quality

To analyse the efficiency on real image quality datasets, we use the following five publicly available datasets with subjective quality ratings, which are often used in image quality prediction research. It should be noted that the datasets come from different sources and each individual dataset may contain unique distortions that are not present in the other datasets.

CSIQ [21]: 866 images in total from 30 reference images, each distorted using one of five types of distortions. **TID 2013** [22]: 3,750 images in total from 24 reference images and 25 types of distortions. **Live Challenge** [23]: 1,162 images with authentic image distortions captured using a representative variety of modern mobile devices. **Live IQA R2** [24]: 779 images in total with 5 different distortion types. **Live MD** [25]: 450 images in total, containing two types of multiply distorted images.

We again analysed the proposed algorithm with a leave-one-dataset-out cross-validation, where we trained the model on the four remaining datasets. As the datasets were rated on different scales, we linearly re-scaled all ratings to a range from 1–5, where 5 is the highest possible quality. Then, for each training run, we saved the model of the epoch on which the best results in terms of PCC was achieved on the validation dataset. As prediction model, we used Pytorch’s ResNet50 implementation with 50 layers, where we replaced the last fully connected

layer with a fully connected layer with only one output. Due to the large variety of different distortion in the datasets, we used the ImageNet pretrained weights and then fine-tuned the entire model during training. Because ResNet expects images of size 224x224x3, we only considered the center crop of the validation images to simplify the evaluation (instead of, e.g., averaging over multiple crops). However, to improve the accuracy of the predictions, a random crop on the training data was applied to increase the training data variety.

We ran the training 15 times with an early stop of 20 epochs on the validation PCC, a learning rate of 0.0001, and a mini-batch size of 32. The average results in terms of PCC are presented in Table I. On four of the five datasets, the proposed loss outperforms the vanilla MSE loss function, whereas, for dataset TID 2013, the vanilla MSE achieves the same results. The correlation of the Live Challenge dataset is overall very low with a PCC of 0.49/0.50. The low performance can be explained by the different types of distortions in this dataset (e.g., lens flare) compared to the training datasets. On the CSIQ dataset, a performance increase of 0.02 in terms of PCC can be observed, showing that the image quality prediction can be improved by applying the proposed bias-aware loss.

TABLE I: Validation results as average PCC of 15 training runs.

Validation Dataset	MSE-Loss	Proposed
Synthetic	0.77	0.93
Speech Quality: EXP1a	0.95	0.96
Speech Quality: EXP1d	0.95	0.97
Speech Quality: EXP1o	0.96	0.95
Speech Quality: EXP3a	0.87	0.92
Speech Quality: EXP3c	0.89	0.90
Speech Quality: EXP3d	0.89	0.92
Image Quality: CSIQ	0.82	0.85
Image Quality: Live Challenge	0.50	0.52
Image Quality: Live IQA	0.88	0.89
Image Quality: Live MD	0.81	0.83
Image Quality: TID 2013	0.70	0.70

IV. CONCLUSION

We presented an open-sourced loss function that automatically learns biases that occur when subjective quality experiments are conducted. The bias-aware loss does not punish prediction errors caused by biases, while the rank order within the dataset is predicted correctly. We could show, on the basis of a synthesised dataset, that the proposed method obtains almost the same results as if the datasets were not exposed to any biases and outperformed the vanilla MSE loss by a relative improvement of 21%. The performance on real data largely depends on the datasets that it is applied to. In cases where no biases are present between datasets, the model gives similar results as a vanilla MSE loss, however, in most cases the results could be improved. Because the bias-aware loss can be applied without additional data or computational costs, it is a helpful tool to improve speech, image, or video quality prediction models that are trained from multiple datasets.

REFERENCES

- [1] W. Lin and C.-C. J. Kuo, "Perceptual visual quality metrics: A survey," *Journal of visual communication and image representation*, vol. 22, no. 4, pp. 297–312, 2011.
- [2] A. Mittal, R. Soundararajan, and A. C. Bovik, "Making a "completely blind" image quality analyzer," *IEEE Signal Process. Lett.*, vol. 20, no. 3, pp. 209–212, 2012.
- [3] A. Mittal, A. K. Moorthy, and A. C. Bovik, "No-reference image quality assessment in the spatial domain," *IEEE Trans. Image Process.*, vol. 21, no. 12, pp. 4695–4708, 2012.
- [4] Zhou Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Trans. Image Process.*, vol. 13, no. 4, pp. 600–612, 2004.
- [5] G. Mittag, B. Naderi, A. Chehadi, and S. Möller, "NISQA: A deep CNN-self-attention model for multidimensional speech quality prediction with crowdsourced datasets," 2021.
- [6] A. A. Catellier and S. D. Voran, "Wawenets: A no-reference convolutional waveform-based approach to estimating narrowband and wideband speech quality," in *ICASSP 2020*, 2020.
- [7] X. Dong and D. S. Williamson, "An attention enhanced multi-task model for objective speech assessment in real-world environments," in *ICASSP 2020*, 2020.
- [8] C. Lo, S. Fu, W.-C. Huang, X. Wang, J. Yamagishi, Y. Tsao, and H.-M. Wang, "MOSNet: Deep Learning-Based Objective Assessment for Voice Conversion," in *Interspeech 2019*, 2019.
- [9] G. Mittag and S. Möller, "Deep Learning Based Assessment of Synthetic Speech Naturalness," in *Interspeech 2020*, 2020.
- [10] S. Bosse, D. Maniry, K.-R. Müller, T. Wiegand, and W. Samek, "Deep neural networks for no-reference and full-reference image quality assessment," *IEEE Trans. Image Process.*, vol. 27, no. 1, pp. 206–219, 2017.
- [11] D. Chen, Y. Wang, and W. Gao, "No-reference image quality assessment: An attention driven approach," *IEEE Trans. Image Process.*, 2020.
- [12] H. Ren, D. Chen, and Y. Wang, "RAN4IQA: restorative adversarial nets for no-reference image quality assessment," in *AAAI 2018*, 2018.
- [13] B. Naderi and S. Möller, "Transformation of mean opinion scores to avoid misleading of ranked based statistical techniques," in *QoMEX 2020*, 2020.
- [14] ITU-T Rec. P.1401, "Methods, metrics and procedures for statistical evaluation, qualification and comparison of objective quality prediction models," 2001.
- [15] S. Zielinski, F. Rumsey, and S. Bech, "on some biases encountered in modern audio quality listening tests-a review," *journal of the audio engineering society*, vol. 56, no. 6, pp. 427–451, june 2008.
- [16] S. Möller, A. Raake, N. Kitawaki, A. Takahashi, and M. Waltermann, "Impairment factor framework for wide-band speech codecs," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 14, no. 6, pp. 1969–1976, 2006.
- [17] P. Kabal, "TSP speech database," McGill University, Quebec, Canada, Tech. Rep. Database Version 1.0, 2002.
- [18] ITU-T Rec. G.107, "The E-model: A computational model for use in transmission planning," 2006.
- [19] G. Mittag and S. Möller, "Non-intrusive speech quality assessment for super-wideband speech communication networks," in *ICASSP 2019*, 2019.
- [20] ITU-T P Suppl. 23, "ITU-T coded-speech database," 1998.
- [21] E. C. Larson and D. M. Chandler, "Most apparent distortion: full-reference image quality assessment and the role of strategy," *Journal of electronic imaging*, vol. 19, no. 1, 2010.
- [22] N. Ponomarenko, L. Jin, O. Ieremeiev, V. Lukin, K. Egiazarian, J. Astola, B. Vozel, K. Chehdi, M. Carli, F. Battisti, and C.-C. J. Kuo, "Image database TID2013: Peculiarities, results and perspectives," *Signal Processing: Image Communication*, vol. 30, pp. 57 – 77, 2015.
- [23] D. Ghadiyaram and A. C. Bovik, "Massive online crowdsourced study of subjective and objective picture quality," *IEEE Trans. Image Process.*, vol. 25, no. 1, pp. 372–387, 2016.
- [24] L. C. H.R. Sheikh, Z.Wang and A. Bovik, "Live image quality assessment database release 2," 2005. [Online]. Available: <http://live.ece.utexas.edu/research/quality>
- [25] D. Jayaraman, A. Mittal, A. K. Moorthy, and A. C. Bovik, "Objective quality assessment of multiply distorted images," in *ASILOMAR 2012*, 2012, pp. 1693–1697.