# AlignNet: Learning dataset score alignment functions to enable better training of speech quality estimators

Jaden Pieper<sup>1</sup>, Stephen Voran<sup>1</sup>

<sup>1</sup>Institute for Telecommunication Sciences, National Telecommunications and Information Administration, United States

jpieper, svoran@ntia.gov

#### **Abstract**

We develop two complementary advances for training noreference (NR) speech quality estimators with independent datasets. Multi-dataset finetuning (MDF) pretrains an NR estimator on a single dataset and then finetunes it on multiple datasets at once, including the dataset used for pretraining. AlignNet uses an AudioNet to generate intermediate score estimates before using the Aligner to map intermediate estimates to the appropriate score range. AlignNet is agnostic to the choice of AudioNet so any successful NR speech quality estimator can benefit from its Aligner. The methods can be used in tandem, and we use two studies to show that they improve on current solutions: one study uses nine smaller datasets and the other uses four larger datasets. AlignNet with MDF improves on other solutions because it efficiently and effectively removes misalignments that impair the learning process, and thus enables successful training with larger amounts of more diverse data.

**Index Terms**: corpus effect, listening experiment, machine learning, no-reference estimator, speech quality, subjective test

#### 1. Background

Speech quality, naturalness, and related quantities can be measured in listening experiments or estimated by algorithms. Algorithms can save much time and effort, but it is a significant challenge to develop an algorithm that produces reliable estimates across a wide range of conditions [1]. The most useful algorithms are called "no-reference" (NR) because they measure impaired speech directly and do not need a reference speech signal for comparison. NR algorithms parallel absolute category rating (ACR) listening experiments where listeners score impaired speech without comparing to reference speech.

NR estimators originally used explicit models [2–5] but moved to data-driven implicit modeling as machine learning (ML) became more mature and practical. A few examples include [6–12]. The ML approach is powerful and effective, but also highly dependent on the quantity and diversity of listening experiment results that comprise the ground-truth training data. This motivates us to combine the results of multiple listening experiments to achieve the needed quantity and diversity.

#### 1.1. The alignment problem

Despite the name, ACR results are not truly absolute. They can depend on a variety of factors including characteristics of individual listeners and the range of conditions included in an experiment [13–16]. This is sometimes called the "corpus effect." For example, in [15] a single fixed synthesized voice file received quality scores ranging from 1.8 to 4.5 in five different experiments, where the only difference between those experiments was the range of conditions included in each. More generally, for one of the five experiments in [15] 41% of the speech

files move by more than 1.0 on the five-point MOS scale when rated in one of the other four experiments. This behavior can be attributed to the listeners' desire to use the entire scale.

The corpus effect creates a problem when we seek to combine results of multiple experiments. Naive combining can yield inconsistent training data that harms training processes instead of enhancing them. But if the results of each experiment can be brought to a single common scale ("aligned"), they can then work together to improve training of an NR estimator. This requires finding a useful common scale and an optimal mapping from each set of experiment results to that common scale. This is the dataset alignment problem.

#### 1.2. Prior work and our contributions

Several alignment techniques have been proposed and used. When two listening experiments include common conditions, results for these conditions may be used to develop alignment functions. Historically many experiments included a standardized adjustable reference condition called the modulated noise reference unit (MNRU) [17] for exactly this purpose. This is currently not common practice, likely because the impairments produced by the MNRU sound very different from impairments appearing in current experiments, thus limiting its usefulness as a reference condition. Other standardized reference conditions have been used in the past but devoting experiment conditions to references always consumes precious resources.

An iterative approach that alternately optimizes alignment functions and an estimation algorithm is given in [18]. An updated iterative approach that leverages ML is given in [19]. Other approaches explore individual alignments for each listener in a listening test [20–22]. Aligning listeners can compensate for their individual behaviors and can lead to better training of NR estimators. But this cannot account for the primary portion of the corpus effect, caused by the broader biases due to each experiment's context. Further, listener alignment is not possible unless datasets label each listener's scores.

In this paper we offer the following novel contributions:

- multi-dataset finetuning, a progressive training regimen that advantageously leverages both larger and smaller datasets
- adding a small score alignment network and a dataset indicator to an audio network
- combining these to learn embeddings for the dataset indicator, alignment functions, and optimal audio network weights
- using an unprecedented 13 datasets covering 3 languages, scores for 4 different speech attributes, and a very wide range of measurement domains, totalling over 300 hours of speech
- demonstrating that these innovations allow previously incompatible datasets to collaborate during training, resulting in better estimates across disparate measurement domains.

# 2. No-reference speech quality estimators for multiple datasets

Here we discuss issues with the conventional strategy for training a speech quality estimator using multiple datasets. We then propose two innovations that enable learning meaningful relationships between audio and scores across multiple datasets, even when inconsistent scores are present.

#### 2.1. Conventional approach

The conventional approach to training an NR speech quality estimator with speech and scores from multiple distinct listening experiments is to simply use all the datasets at once. However, due to the corpus effect, there can be a misalignment between speech and target scores from different experiments. When identical or very similar speech files appear in multiple listening experiments, they almost certainly receive different scores in each experiment. This means that while training, the network must attempt to map identical or similar input files to multiple conflicting output scores, and would likely estimate a score that is roughly the average of all scores seen for the file. This is reasonable to an extent, but these disparate scores for the same input add additional noise for the network to sift through and place inherent limitations on its estimates; it can never produce a single estimate for this input that achieves low loss for all the associated target scores.

#### 2.2. Multi-dataset finetuning

In the conventional approach the network attempts to learn audio relationships while dealing with misaligned target scores from multiple experiments, which impedes the training process. Previous work has demonstrated the benefits of pretraining a network on one dataset and then switching to a different dataset for finetuning [1, 23, 24]. The initial pretraining allows the network to learn a basic and somewhat transferable relationship between audio and scores. Here we propose multi-dataset finetuning (MDF), where we first pretrain the network on a single dataset before finetuning with all the datasets at once, including the original dataset used for pretraining. Pretraining places the network into a state where it already knows some meaningful relationships between audio and scores, and can then balance the misaligned scores from the different listening experiments. It has some of the same limitations as the conventional approach, but pretraining on a single dataset enables much better training and predictions when applied to multiple datasets. We believe MDF is a novel approach to the problem.

#### 2.3. Dataset alignment with AlignNet

We now introduce a novel architecture called AlignNet, which allows any NR speech quality estimator to better benefit from multiple datasets, with only a minimal increase in network complexity. AlignNet is essentially two network components in sequence, which we call the AudioNet and the Aligner respectively. AlignNet is intentionally designed to be agnostic to the choice of the AudioNet, which maps audio or audio features (depending on the choice of AudioNet) to intermediate score estimates. The Aligner uses a categorical dataset indicator to map those intermediate score estimates to final scores for the appropriate dataset. The network architecture is outlined in Fig. 1.

It is necessary to select a reference dataset and force the Aligner to apply the identity function to reference dataset results. This ensures that the outputs of AudioNet are grounded in meaningful quality scores, and that AudioNet gives estimates for all audio in the domain of the reference dataset scores. Favorable attributes for a reference dataset include trustworthy

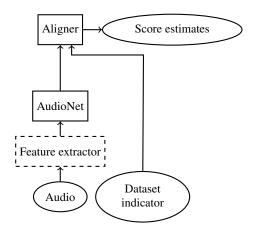


Figure 1: AlignNet model diagram. The feature extractor is optional depending on the choice of AudioNet.

scores and a wide range of conditions. This allows AudioNet to produce outputs on a level playing field for comparisons of speech files, without any corpus effect biases muddling the rankings. Finally this ensures that the mappings learned by the Aligner carry relevant information about the relationships between different experiments and are easily interpretable.

To build an intuition, again consider multiple listening experiments where very different scores were reported for an identical condition and audio file. In AlignNet, the AudioNet would give the same result for each occurrence of the audio file, and then the Aligner could use the dataset indicator to successfully map this single value to the appropriate, different scores. Thus AlignNet is able to yield low loss for each occurrence of the input, in spite of their different scores.

The Aligner is extremely light-weight and adds an insignificant number of parameters to any effective AudioNet. It first maps the dataset indicator to an N-dimensional embedding and concatenates that embedding with the AudioNet estimation. The network consists of a series of fully connected layers of identical dimensions, separated by ReLU activations, and a final fully connected layer that maps the data to a single score estimate. The number of parameters in the Aligner is dependent on the number of datasets used, and in our implementations it has roughly 1100 parameters total. More in-depth implementation details are available with our source code  $^1$ .

Successful training of AlignNet requires a clear division of labor — the AudioNet should not attempt to make alignments and the Aligner should not attempt to measure audio. This is achieved by MDF, the use of a reference dataset, and freezing the AudioNet for one epoch at the start of finetuning. The reference dataset is used for pretraining so AudioNet is well-positioned at the beginning of training, at least for that dataset. Continuing to train allows the AudioNet to better learn the speech from the other datasets, while the Aligner learns how to reconcile the scores across the different datasets.

#### 2.4. Loss function

In conventional network training each individual piece of data is given equal weight in the loss function. When training a network with multiple datasets, we propose that each *dataset* be given equal weight within the loss function using

$$L = \frac{1}{N_d} \sum_{j=1}^{N_d} l(\mathbf{y}_j, \hat{\mathbf{y}}_j), \tag{1}$$

https://github.com/NTIA/alignnet

Table 1: Summary of 13 datasets used in this work. The first 9 have 8 to 24 votes per file, the final 4 have 4 or 5 votes per file on average. Blizzard 2021 uses Spanish language, Tencent uses Chinese, all others use English.

Dataset	Abbr.	Domain	Number of Files	
Blizzard 2021 SS1 [25]	B21 S1 Nat	Synthesized & natural speech	242	
Blizzard 2021 SH1 [25]	B21 H1 Nat	same as above	338	
Blizzard 2021 SS1 [25]	B21 S1 Acc	same as above	363	
Blizzard 2008 News [26]	B08 News	same as above	802	
Blizzard 2008 Novels [26]	B08 Novel	same as above	802	
FFTnet [27]	FFTnet	Neural vocoders	1200	
NOIZEUS [28]	NOIZEUS	Noise & suppression	1664	
VoiceMOS Challenge 2022 [29]	VMC22	Speaker conversion & synthesized speech	7106	
Tencent [30]	Tencent	Noise, suppression, reverb, coding, packet loss & concealment	11,563	
NISQA SIM [31]	NISQA	Coding, packet loss, noise, filtering & clipping	12,500	
Voice Conversion Challenge 2018 [32]	VCC18	Voice conversion systems	20,580	
Indiana U. MOS [33]	IU MOS	Noise & reverb	36,000	
PSTN [34]	PSTN	PSTN to VoIP calls plus noise	58,709	

where  $N_d$  is the number of datasets used in training,  $\mathbf{y}_j$  are all the targets for dataset j,  $\hat{\mathbf{y}}_j$  are all the estimates for dataset j, and l is mean-squared error loss. This allows a network to achieve good results for each dataset, rather than letting larger datasets dominate the learning.

## 3. Experiments and Results

We performed studies with two different groups of datasets: one with nine smaller datasets that have more votes per file and one with four larger datasets that have fewer votes per file. Key properties are summarized in Table 1. Three languages are represented and the total duration of speech exceeds 300 hours. Diverse measurement domains include synthesized speech, voice conversion, neural vocoders, conventional codecs, packet loss, noise, reverb, enhancement, filtering, and more. The NOIZEUS and PSTN datasets are narrowband (nominally 300 to 3600 Hz) and the remaining datasets support wideband (nominally 50 to 7000 Hz) or fullband speech. Scores are for four different attributes: "Acceptability" (B21 S1 Acc), "Naturalness" (remaining 4 Blizzard datasets, VMC22, and VCC18), "Overall Quality" (NOIZEUS), and "Speech Quality" (remaining 5 datasets). These attributes are related but not identical which further motivates dataset alignment.

In each study we explored different training regimens and network architectures with multiple datasets in order to demonstrate the performance of our two novel approaches compared to existing methods. We used the NR speech quality estimator MOSNet [35] for all experiments in both studies, either exclusively or as the AudioNet in AlignNet. We chose MOSNet as it is a sufficiently large network to have the capacity to achieve good results for these studies, while being small enough to train relatively quickly and be more usable in practice. Further it has been successfully used as a baseline model in related research on listener specific corrections [21, 22]. Unlike the original MOSNet implementation, we did not opt to use frame level loss, and instead averaged frames into a single value before the loss function; otherwise our implementation exactly matches the original paper. All audio was resampled to 16 kHz prior to the STFT calculation. We randomly split each dataset into 80%, 10%, and 10% for the training, validation, and testing data respectively. We use the same split across all tests to ensure a fair comparison, and all reported results are from the unseen test sets. We trained all networks with the loss function defined in (1), except in the bias-aware loss comparison (BAL), which uses the loss function defined in [19]. In the small dataset study the Tencent dataset was selected as the reference dataset, which means we also used it for MDF pretraining. In the large dataset study the NISQA dataset filled this role.

We use "depth" to describe network performance for a single dataset and "breadth" to describe performance across all datasets of interest. Each AudioNet has its natural tradeoff between depth and breadth; at a certain point one cannot be improved without harming the other. Adding dataset alignment can mitigate this tradeoff and allow better simultaneous depth and breadth. AlignNet's Aligner generally improves depths without reducing breadth.

We use two metrics to evaluate the depth and breadth performance in our experiments: Pearson's linear correlation coefficient (LCC) describes the networks ability to rank speech attributes, and root mean-squared error (RMSE) describes the distance of the network's estimates from the true scores. The small and large dataset study results are given in Tables 2 and 3, respectively. Each column gives results for the unseen testing portion of a given dataset. Bold indicates the best and underlining the second best performance for each column. The \* symbol denotes a statistically significant improvement over the conventional regimen ("All" row), calculated using Zou's confidence interval for LCC [36] and bootstrapping for RMSE [37]. For the large dataset study full results are shown for training on each dataset and the diagonal is shaded; off-diagonal results show the lack of breadth. We did the same for the small dataset study but for brevity the diagonal is compressed into the "Individual" row, where each cell shows results for training and testing (on unseen data) for a given dataset. No individually trained model shows any meaningful breadth and none of the unshown values were first or second place for any column. Note that it is more difficult to achieve statistical significance for the datasets with fewer than 500 files, as the test sets are very small.

### 3.1. Training regimens and network architectures

We demonstrate the improvements provided by MDF and AlignNet through comparison with a series of baselines. The first set of baselines trained MOSNet on each target dataset individually. We also trained MOSNet with the conventional approach discussed in Sec. 2.1, which is denoted as "All" in the results tables. The MDF approach is denoted as "All (+ MDF)."

We also trained MOSNet using the bias-aware loss (BAL) method defined in [19]. BAL seeks to address the corpus effect in training by using least-squares to estimate a scale and shift for each dataset after each training epoch. The scale and shift are then used in the loss function to attempt to harmonize the disparate experiment scores. BAL training relies on a hyperparameter  $r_{\rm th}$ . Scale and shift are not used in the loss calculation until training correlation exceeds  $r_{\rm th}$ . We selected a single threshold of  $r_{\rm th}=0.6$  based on curves shown in [19]. We use the same dataset as reference for AlignNet and BAL.

We implemented AlignNet with MOSNet as the AudioNet and an Aligner that uses a 10-dimensional dataset embedding and has 5 fully connected layers of dimension 16. MOSNet has roughly 1.2 million parameters and the Aligner has roughly 1100, meaning the Aligner was less than 0.1% of the total network size. To encourage the Aligner to focus only on dataset alignment we froze the pretrained AudioNet for the first epoch. In the tables AlignNet is denoted as "All (+ MDF + AlignNet)". Neither MDF nor AlignNet add measurable additional training time. BAL increases training time by about 50%.

Table 2: LCC (above) and RMSE (below) for all models on the small datasets.

Training Data	B21 S1 Nat	B21 H1 Nat	B21 S1 Acc	B08 Novel	B08 News	FFTNet	NOIZEUS	VMC22	Tencent	All
Individual	0.45	0.65	0.23	0.66	0.67	0.81*	0.80*	0.51	0.94*	NA
All	0.73	0.82	0.70	0.62	0.69	0.53	0.65	0.71	0.80	0.77
All (+ BAL)	0.83	0.77	0.65	0.56	0.55	0.69*	0.70	0.75*	0.91*	0.83*
All (+ MDF)	0.82	0.83	0.62	0.72	0.77	0.70*	0.71	0.74	0.89*	0.84*
All (+ MDF + AlignNet)	0.88	0.90	0.78	0.81*	0.82*	0.66*	<u>0.76*</u>	0.76*	0.92*	0.87*
Training Data	B21 S1 Nat	B21 H1 Nat	B21 S1 Acc	B08 Novel	B08 News	FFTNet	NOIZEUS	VMC22	Tencent	All
Training Data Individual	<b>B21 S1 Nat</b> 1.15	<b>B21 H1 Nat</b> 1.19	<b>B21 S1 Acc</b> 0.95	<b>B08 Novel</b> 0.88	<b>B08 News</b> 0.87	<b>FFTNet</b> 0.66*	NOIZEUS 0.34*	VMC22 0.95	Tencent 0.41*	All NA
Individual	1.15	1.19	0.95	0.88	0.87	0.66*	0.34*	0.95	0.41*	NA
Individual All	1.15 0.87	1.19 <u>0.69</u>	0.95 <u>0.56</u>	0.88	0.87 0.70	0.66*	0.34* 0.44	0.95 0.66	<b>0.41*</b> 0.80	NA 0.73

Table 3: LCC (above) and RMSE (below) for all models on the large datasets.

Training Data	NISQA	VCC18	IU	PSTN	All
NISQA	0.89*	0.48	0.69	0.74	0.39
VCC18	0.50	0.66*	0.42	0.53	0.07
IU	0.54	0.13	0.97*	0.58	0.50
PSTN	0.72	0.38	0.62	0.81*	0.37
All	0.80	0.60	0.86	0.75	0.88
All (+ BAL)	0.88*	0.62	0.96*	0.80*	0.94*
All (+ MDF)	0.91*	0.63*	0.96*	0.81*	0.94*
All (+ MDF + AlignNet)	0.91*	0.64*	0.97*	0.80*	0.94*
Training Data	NIICOA	X70040	TTT	DOTAL	4 11
Training Data	NISQA	VCC18	IU	PSTN	All
NISQA	0.53*	1.15	3.50	0.92	2.02
NISQA	0.53*	1.15	3.50	0.92	2.02
NISQA VCC18	0.53* 0.98	1.15 <b>0.71</b> *	3.50 4.19	0.92	2.02
NISQA VCC18 IU	0.53* 0.98 3.46	1.15 0.71* 3.52	3.50 4.19 <b>0.48</b> *	0.92 0.99 2.65	2.02 2.36 2.54
NISQA VCC18 IU PSTN	0.53* 0.98 3.46 0.82	1.15 0.71* 3.52 1.03	3.50 4.19 <b>0.48*</b> 3.00	0.92 0.99 2.65 <b>0.51</b> *	2.02 2.36 2.54 1.70
NISQA VCC18 IU PSTN All	0.53* 0.98 3.46 0.82	1.15 0.71* 3.52 1.03 0.85	3.50 4.19 <b>0.48*</b> 3.00 1.14	0.92 0.99 2.65 <b>0.51*</b> 0.92	2.02 2.36 2.54 1.70 0.97

#### 3.2. Results

As expected, training on individual datasets gives reasonable depth for some of the larger datasets, but almost no breadth. Further, only datasets with over 1000 files demonstrate good depth, and some of the particularly small datasets fail to train meaningfully at all. The conventional approach ("All" row) offers an improvement over individual training for the smaller datasets — it gives better breadth as one would expect. However, outside of VMC22, for datasets with more than 1000 files "All" is significantly worse than individual training. Individual training is particularly successful for FFTNet and NOIZEUS because those datasets contain very unique conditions.

MDF has no specific features to address dataset alignment but it often outperforms or matches the BAL method, including when RMSE and correlations are measured across all the data in the small dataset study ("All" column in Table 2). The combination of pretraining and a loss function that gives equal weight to each dataset enables this network to reliably estimate scores with no per-dataset side information.

AlignNet demonstrates the best performance of all the training regimens and model architectures that are trained on all the datasets, particularly in the small dataset study. There it achieved the highest LCCs for every dataset but two and the lowest RMSE for every dataset but one. It also achieved the best performance in the "All" column for both metrics in both studies. This demonstrates that AlignNet is a powerful tool for reconciling different datasets. Remarkably, when compared to training on individual datasets, AlignNet gives better or similar results on *all datasets at once*. AlignNet is an estimator that is capable of properly ranking a multitude of diverse audio conditions, with an insignificant increase to overall model

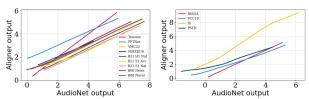


Figure 2: Learned dataset score alignment functions. Left-small dataset study. Right - large dataset study. Functions plotted only over observed values in training data for each dataset. complexity compared to only using the AudioNet. In addition to better estimates, AlignNet offers a few other benefits over BAL. AlignNet simultaneously updates the AudioNet and the Aligner, which is far more efficient than the iterative approach of BAL. Further, [19] states that BAL is very sensitive to the  $r_{\rm th}$  parameter and requires optimization for each dataset, which becomes impractical with a large total number of files. AlignNet has a similar hyperparameter which sets the duration that AudioNet is frozen once MDF starts, but we consistently see best performance with the fixed value of one epoch.

It is easy to visualize the learned alignments from AlignNet by plotting the Aligner score estimates vs intermediate score estimates from AudioNet, as seen in Fig. 2. Datasets with similar properties give similar alignment functions, which can be seen clearly in the alignment function plot for the small dataset study. After they have been learned, these alignments could be approximated by monotonic third-degree polynomials, as previously recommended by [14]. Finally, note that there can be additional information carried in the range of the intermediate scores of the non-reference datasets. These scores can extend beyond the nominal range which may speak to the impairments in those datasets relative to those in the reference dataset.

#### 4. Conclusion

AlignNet with MDF can reconcile different rated speech attributes such as naturalness, acceptability, and quality. In the small dataset study four attributes were successfully harmonized resulting in a more robust NR estimator and revealing the relationships between the different attributes and experiment contexts. This work demonstrates that disparate scores from distinct listening experiments can be used harmoniously for NR speech estimator training by adding a small alignment network to an existing NR speech estimator. This type of work is always limited by the data used. We argue our work is very strong in this regard, but we nonetheless seek to do additional work with even more and broader data. We also plan to experiment with other choices for the AudioNet used inside of AlignNet and to further interpret the learned relationships between datasets. Finally, studying the performance of an AlignNet model with MDF on unseen datasets could provide additional insights into the practical use of such a model.

#### 5. References

- [1] E. Cooper, W. C. Huang, T. Toda, and J. Yamagishi, "Generalization ability of MOS prediction networks," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing*, 2021, pp. 8442–8446.
- [2] J. Liang and R. Kubichek, "Output-based objective speech quality," in *Proc. IEEE Vehicular Technology Conference*, vol. 3, June 1994, pp. 1719–1723.
- [3] D. S. Kim, "ANIQUE: An auditory model for single-ended speech quality estimation," *IEEE Trans. on Speech and Audio Processing*, vol. 13, no. 5, pp. 821–831, 2005.
- [4] L. Malfait, J. Berger, and M. Kastner, "P.563—The ITU-T standard for single-ended speech quality assessment," *IEEE Trans. on Audio, Speech, and Language Processing*, vol. 14, no. 6, pp. 1924–1934, Nov. 2006.
- [5] T. H. Falk, C. Zheng, and W.-Y. Chan, "A non-intrusive quality and intelligibility measure of reverberant and dereverberated speech," *IEEE Trans. on Audio, Speech, and Language Processing*, vol. 18, no. 7, pp. 1766–1774, Sept. 2010.
- [6] M. H. Soni and H. A. Patil, "Novel deep autoencoder features for non-intrusive speech quality assessment," in *Proc. European Signal Processing Conference*, Nov. 2016, pp. 2315–2319.
- [7] M. Hakami and W. B. Kleijn, "Machine learning based nonintrusive quality estimation with an augmented feature set," in Proc. IEEE International Conference on Acoustics, Speech and Signal Processing, 2017, pp. 5105–5109.
- [8] S. Fu, Y. Tsao, H. Hwang, and H. Wang, "Quality-Net: An end-to-end non-intrusive speech quality assessment model based on BLSTM," in *Proc. Interspeech*, 2018, pp. 1873–1877.
- [9] H. Gamper, C. K. A. Reddy, R. Cutler, I. J. Tashev, and J. Gehrke, "Intrusive and non-intrusive perceptual speech quality assessment using a convolutional neural network," in *Proc. IEEE Workshop* on Applications of Signal Processing to Audio and Acoustics, 2019, pp. 85–89.
- [10] C. K. A. Reddy, V. Gopal, and R. Cutler, "DNSMOS P.835: A non-intrusive perceptual objective speech quality metric to evaluate noise suppressors," in *Proc. IEEE International Conference* on Acoustics, Speech and Signal Processing, 2022, pp. 886–890.
- [11] R. E. Zezario, S.-W. Fu, F. Chen, C.-S. Fuh, H.-M. Wang, and Y. Tsao, "Deep learning-based non-intrusive multi-objective speech assessment model with cross-domain features," *IEEE/ACM Trans. on Audio, Speech, and Language Processing*, vol. 31, pp. 54–70, 2023.
- [12] A. A. Catellier and S. D. Voran, "Wideband audio waveform evaluation networks: Efficient, accurate estimation of speech qualities," *IEEE Access*, vol. 11, pp. 125 576–125 592, 2023.
- [13] ITU-T Recommendation P.800.2, Mean opinion score interpretation and reporting, Geneva, 2016.
- [14] ITU-T Recommendation P.1401, Methods, metrics and procedures for statistical evaluation, qualification and comparison of objective quality prediction models, Geneva, 2001.
- [15] E. Cooper and J. Yamagishi, "Investigating range-equalizing bias in mean opinion score ratings of synthesized speech," in *Proc. Interspeech*, 2023, pp. 1104–1108.
- [16] S. Le Maguer, S. King, and N. Harte, "Back to the future: Extending the Blizzard Challenge 2013," in *Proc. Interspeech*, 2022, pp. 2378–2382.
- [17] ITU-T Recommendation P.810, "Modulated noise reference unit (MNRU)," Geneva, 1996.
- [18] S. Voran, "An iterated nested least-squares algorithm for fitting multiple data sets," U.S. Department of Commerce, National Telecommunications and Information Administration, Institute for Telecommunication Sciences, Tech. Rep. TM-03-397, 2002.
- [19] G. Mittag, S. Zadtootaghaj, T. Michael, B. Naderi, and S. Möller, "Bias-aware loss for training image and speech quality prediction models from multiple datasets," in *Proc. Thirteenth International Conference on Quality of Multimedia Experience*, 2021, pp. 97– 102.

- [20] N. Nessler, M. Cernak, P. Prandoni, and P. Mainar, "Non-intrusive speech quality assessment with transfer learning and subjectspecific scaling," in *Proc. Interspeech*, 2021, pp. 2406–2410.
- [21] Y. Leng, X. Tan, S. Zhao, F. Soong, X.-Y. Li, and T. Qin, "MB-NET: MOS prediction for synthesized speech with mean-bias network," in *Proc. IEEE International Conference on Acoustics*, Speech and Signal Processing, 2021, pp. 391–395.
- [22] W.-C. Huang, E. Cooper, J. Yamagishi, and T. Toda, "LDNet: Unified listener dependent modeling in MOS prediction for synthetic speech," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing*, 2022, pp. 896–900.
- [23] W. C. Tseng, C. Y. Huang, W. T. Kao, Y. Y. Lin, and H. Y. Lee, "Utilizing self-supervised representations for MOS prediction," in *Proc. Interspeech*, 2021, pp. 3521–3525.
- [24] H. Becerra, A. Ragano, and A. Hines, "Exploring the influence of fine-tuning data on wav2vec 2.0 model for blind speech quality prediction," in *Proc. Interspeech*, 2022, pp. 4088–4092.
- [25] Z.-H. Ling, X. Zhou, and S. King, "The Blizzard challenge 2021," in *Proc. Blizzard Challenge Workshop*, 2021.
- [26] V. Karaiskos, S. King, R. A. J. Clark, and C. Mayo, "The Blizzard challenge 2008," in *Proc. Blizzard Challenge Workshop*, 2008.
- [27] Z. Jin, A. Finkelstein, G. J. Mysore, and J. Lu, "FFTNet: a realtime speaker-dependent neural vocoder," in *Proc. IEEE Interna*tional Conference on Acoustics, Speech and Signal Processing, 2018
- [28] Y. Hu and P. Loizou, "Subjective comparison of speech enhancement algorithms," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing*, 2006.
- [29] W. C. Huang, E. Cooper, Y. Tsao, H.-M. Wang, T. Toda, and J. Yamagishi, "The VoiceMOS Challenge 2022," in *Proc. Interspeech* 2022, 2022, pp. 4536–4540.
- [30] G. Yi, W. Xiao, Y. Xiao, B. Naderi, S. Moller, W. Wardah, G. Mittag, R. Cutler, Z. Zhang, D. S. Williamson, F. Chen, F. Yang, and S. Shang, "ConferencingSpeech 2022 Challenge: Non-intrusive objective speech quality assessment challenge for online conferencing applications," in *Proc. Interspeech*, 2022, pp. 3308–3312.
- [31] G. Mittag, B. Naderi, A. Chehadi, and S. Möller, "NISQA: A deep CNN-self-attention model for multidimensional speech quality prediction with crowdsourced datasets," in *Proc. Inter*speech, 2021, pp. 2127–2131.
- [32] J. Lorenzo-Trueba, J. Yamagishi, T. Toda, D. Saito, F. Villavicencio, T. Kinnunen, and Z. Ling, "The voice conversion challenge 2018: Promoting development of parallel and nonparallel methods," in *Proc. Speaker Odyssey*, 2018.
- [33] X. Dong and D. S. Williamson, "A pyramid recurrent network for predicting crowdsourced speech-quality ratings of real-world signals," in *Proc. Interspeech*, 2020.
- [34] G. Mittag, R. Cutler, Y. Hosseinkashi, M. Revow, S. Srinivasan, N. Chande, and R. Aichner, "DNN no-reference PSTN speech quality prediction," in *Proc. Interspeech*, 2020.
- [35] C.-C. Lo, S.-W. Fu, W.-C. Huang, X. Wang, J. Yamagishi, Y. Tsao, and H.-M. Wang, "MOSNet: Deep learning-based objective assessment for voice conversion," in *Proc. Interspeech*, 2019.
- [36] G. Y. Zou, "Toward using confidence intervals to compare correlations." *Psychological methods*, vol. 12, no. 4, p. 399, 2007.
- [37] B. Efron and R. J. Tibshirani, An introduction to the bootstrap. Chapman and Hall/CRC, 1994.