

BEYOND MOS: SUBJECTIVE IMAGE QUALITY SCORE PREPROCESSING METHOD BASED ON PERCEPTUAL SIMILARITY

Lei Wang, Desen Yuan

University of Electronic Science and Technology of China

ABSTRACT

Image quality assessment often relies on raw opinion scores provided by subjects in subjective experiments, which can be noisy and unreliable. To address this issue, postprocessing procedures such as ITU-R BT.500, ITU-T P.910, and ITU-T P.913 have been standardized to clean up the original opinion scores. These methods use annotator-based statistical priors, but they do not take into account extensive information about the image itself, which limits their performance in less annotated scenarios. Generally speaking, image quality datasets usually contain similar scenes or distortions, and it is inevitable for subjects to compare images to score a reasonable score when scoring. Therefore, In this paper, we proposed Subjective Image Quality Score Preprocessing Method perceptual similarity Subjective Preprocessing (PSP), which exploit the perceptual similarity between images to alleviate subjective bias in less annotated scenarios. Specifically, we model subjective scoring as a conditional probability model based on perceptual similarity with previously scored images, called subconscious reference scoring. The reference images are stored by a neighbor dictionary, which is obtained by a normalized vector dot-product based nearest neighbor search of the images' perceptual depth features. Then the preprocessed score is updated by the exponential moving average (EMA) of the subconscious reference scoring, called similarity regularized EMA. Our experiments on multiple datasets (LIVE, TID2013, CID2013) show that this method can effectively remove the bias of the subjective scores. Additionally, Experiments prove that the Preprocessed dataset can improve the performance of downstream IQA tasks very well.

Index Terms— Blind image quality assessment (BIQA), Subjective Image Quality, Subjective Score Preprocessing, Perceptual Similarity.

1. INTRODUCTION

Image quality assessment [1, 2, 3, 4, 5, 6, 7], assessing the quality of an image is a crucial aspect of image processing research. Its goal is to enable computers to automatically analyze the features of an image and determine its strengths and weaknesses, such as the presence of distortion or other defects. The explosion of data has created a need for efficient

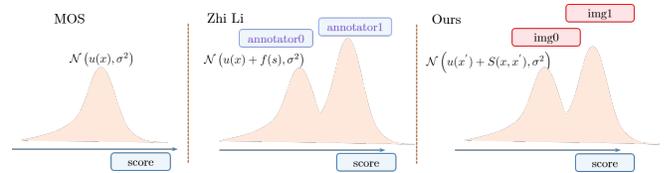


Fig. 1. Left: traditional MOS model, Gaussian Distribution. Middle: model with annotator information added, Gaussian multimodal distribution. Right: model with image information added, Gaussian multimodal distribution. $u(x)$ is regarded as the true quality. $f(s)$ is a learnable bias of annotator s . $u(x')$ is regarded as the true quality of the reference image. $S(x, x')$ represents the score residual converted from the perceptual similarity between images.

processing of large amounts of image and video data. As a result, image quality assessment has become increasingly important for both practical applications and scientific research. It is typically divided into two categories: assessment with reference, and assessment without reference. Reference-based methods require a high-quality reference image to compare against, but this is not always feasible in practice. Therefore, no-reference assessment has become an active area of research in academia and industry.

To reduce the subjectivity of non-reference image quality assessment, this method requires manual labeling of image quality, collection of scores, and data processing to obtain the MOS value. However, there are various errors in human labeling. Currently, the processing methods only take into account the image and personnel numbers and only focus on the statistical results of the data. They ignore the feature information contained in the images themselves. Simply processing the labeling results statistically cannot eliminate the error of personnel labeling, but only provides a smooth process.

The core hypothesis of image quality assessment labeling is that two images with similar image quality have similar labeling scores. However, current MOS processing methods do not pay attention to the hidden information caused by image similarity, which is just the core of the image quality assessment task and does not have task specificity. For the first time, this article introduces image quality correlation to help solve Data postprocessing problems. For the task of image quality

assessment without reference, the objective and unbiased tagging data is the most important for the training of the model.

Such a crowdsourcing way is costly and time-consuming. Some subjective score postprocessing methods have been proposed to replace MOS such as P910 [8], P913 [9], BT500 [10]. However, existing techniques for image quality assessment rely only on the collected scores and ignore the essential information contained in the original images. This can be problematic when the number of annotations is small or when only a single annotator is available. To address this issue, we propose a method named PSP-IQA for postprocessing IQA data based on subconscious reference scoring. It assumes that, When the human visual system makes a subjective quality assessment, it subconsciously gives a score based on a reference image. Specifically, we use the nearest neighbor search (NNS) to find a reference image based on perceptual similarity. Then similarity information in the dataset can help us to modify scores as a regularized item. It takes advantage of the hypothesis that similar image features should have similar scores, thus leveraging the implicit information in the dataset to correct the bias of subjective labeling. Although this regularization term of perceived similarity will limit the fraction away from his similar image. We still believe that the original marked score has a high degree of confidence, so we use EMA to fine-tune the original score by using the perceptual regularization term, which is called perceptual regularization fine-tuning. This simple and effective method is suitable for IQA data postprocessing and is particularly useful when only a limited number of annotations are available. The proposed perceptual similarity postprocessing method is proven to be effective in theory and experiment. We experimented with the IQA task (3 datasets). The results show that the proposed method reduces the adverse effect of subjective bias data on the model performance.

2. RELATED WORK

Image quality assessment.

For image quality assessment datasets [4], there are several acquisition and processing techniques, such as guaranteeing the consistency of the subjective assessment environment [5], adding post-processing to the obtained data [6], and eliminating outliers [7]. Without a doubt, however, several annotators are required for the large-scale subjective assessment datasets that are all amassed through crowdsourcing.

IQA tasks [1, 2, 3], which are frequently thought of as regression problems and have a regressive sigmoid head, have recently been the subject of DNN-based objective assessment algorithms that may perform better. In addition, numerous alternative frameworks for subjective assessment tasks have also been developed using GAN [11], VAE [12], and transformers [13]. For many novel visual tasks and scenarios, such as distorted images [4], virtual reality [14, 15], light fields [16], and dehazing images [17, 18], there is cur-

rently a high demand for fresh datasets. We try to reduce the number of annotations required by the data set by MOS post-processing method, while ensuring the reliability of the model.

Crowdsourced annotations and Noisy label. Any complex function can be learned by an over-parameterized network from corrupted labels [19, 20, 21, 22, 22, 23]. According to Zhang et al. [24], DNNs can easily fit the entire training dataset with any corrupted label ratio, resulting in less generality on the test dataset. Robust loss functions [20, 21, 25], regularization [26], robust network architecture [22], sample selection [23], training strategy [27, 28], and other techniques are proposed to train deep networks in noisy environments.

Most methods, however, are designed around one-hot label properties such as classification task [26] sparsity and noise tolerance [29]. As a result, it cannot be directly applied to the subjective bias problem. Annotators for subjective assessment problems are frequently crowdsourced [30, 31, 32], and each person’s score is biased against the ideal objective evaluator.

Furthermore, the International Telecommunication Union (ITU) and researchers have proposed a number of crowdsourced data processing standards [10, 8, 9] to eliminate subjective bias in MOS. However, these proposed methods do not take into account specific task information but instead rely on iterative fitting to remove subjective bias in the data, the applicability of different tasks is limited, and more data is required, making it difficult to apply under a single annotation.

3. METHODS

In previous papers, the main modeling method of subjective scoring can be expressed as the following formula:

$$p(y|x) = E_{s|x}p(y|x, s)p(s|x) \quad (1)$$

Rather than modeling the scoring $p(y|x)$ as a Gaussian distribution $p(y|x) \sim \mathcal{N}(u(x), \sigma^2)$, additional annotator information is introduced to model the conditional Gaussian distribution $p(y|x, s) \sim \mathcal{N}(u(x) + f(s), \sigma^2)$. It is generally believed that annotators have the same bias for different images. $u(x)$ is regarded as the true quality.

Compared with introducing additional annotator information, we use the perceptual similarity of the image itself as a constraint. We consider the score of the mark to be related to a potential reference image when scoring.

$$p(y|x) = E_{x'|x}p(y|x, x')p(x'|x) \quad (2)$$

where $p(y|x, x') \sim \mathcal{N}(u(x') + S(x, x'), \sigma^2)$, $p(x'|x)$ represents the probability distribution of the potential reference images when the annotator scores. $S(x, x')$ represents the score residual converted from the perceptual similarity between images.

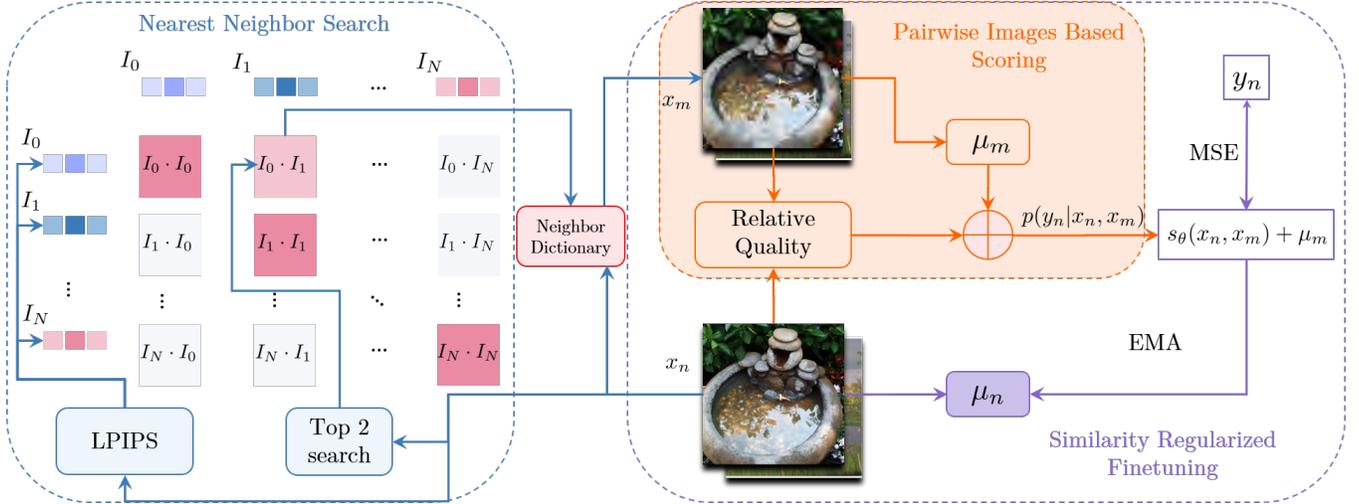


Fig. 2. The overall framework of the proposed method. When the marked score y_n is biased large, the perceptual similarity score and the similar image estimation score can get the correct real score, thus correcting y_n .

In fact, $p(x' | x)$ is difficult to obtain because it is difficult to tell which image the annotator subconsciously compares with. It is natural to think that all images have an equal probability of being adopted by the annotator for reference. However, this would make us computationally unbearable. For simplicity, we model it as a nearest-neighbor search model here. We think people always tend to compare with similar images.

Given a noisy dataset $\{(x_n, y_n)\}_{i=1}^N$, maximizing the log-likelihood can be converted to minimizing the mean squared error (MSE).

$$\arg \min_{\theta} \frac{1}{N} \sum_{n=1}^N (y_n - (S(x_n, \text{NNS}(x_n)) + u(x'_n)))^2, \quad (3)$$

In short, compared to traditional MOS $u(x_n) = y_n$, we added perceptual similarity as a regular term to fine-tune the final score.

Specifically, we use perceptual similarity LPIPS [33] as a metric to search for the most similar images, resulting in NNS. The perceptual similarity score S is obtained by using the features of LPIPS through a scorer, which consists of a ResNet network.

When starting optimization, we initialize $u(x_n)$ as y_n , by defining a score matrix $U \in \mathbb{R}^{1 \times N}$. Iterate the network parameters via gradient descent and update the estimated true score by fine-tuning from y using EMA:

$$\begin{cases} u^{t+1}(x_n) = \text{EMA}(S^t(x_n, x'_n) + u^t(x'_n), u^t(x_n)) \\ \theta^{t+1} = \theta^t - \lambda \Delta \theta \end{cases} \quad (4)$$

where $\Delta \theta$ obtained by backpropagation:

$$\nabla_{\theta} \left\{ \frac{1}{N} \sum_{n=1}^N (y_n - (S^t(x_n, \text{NNS}(x_n)) + u^t(x'_n)))^2 \right\}. \quad (5)$$

In the actual implementation, we set a warm-up time T for the perceptual similarity function S . Before the training epoch times T , we do not adjust the labels. After training the epoch times T times, we use EMA to fine-tune the score:

$$u^{t+1}(x_n) = \begin{cases} \text{EMA}(S^t(x_n, x'_n) + u^t(x'_n), u^t(x_n)) \\ y_n, \text{ when } : t < T \end{cases} \quad (6)$$

4. EXPERIMENTS

In this section, we first describe the experimental setups, including datasets, assessment criteria, and network architecture details. Then we compare the performance of PSP with other preprocessing methods. We next conduct a series of ablation studies to identify the contribution of the key components of PSP. Finally, we also present some visualization cases.

4.1. Datasets and Settings

Datasets. Since there are no existing image quality assessment datasets to measure preprocessing performance with specific annotation information. We selected existing popular quality assessment datasets to generate data with subjective bias, including LIVE [36] and TID2013 [37] and popular quality assessment datasets CID2013 [6] with specific annotation information.

Table 1. We compare the performance of our method with subjective preprocess methods and ablation study on bias label (bias rate 100%) from image quality databases of LIVE, TID2013, and CID2013. We report SROCC, KROCC, and MSE results between the preprocessed quality scores and the true MOS provided by the database. We highlight the best results in bold font.

Datasets	LIVE				TID				CID			
Methods\Metrics	SROCC	PLCC	KROCC	MSE	SROCC	PLCC	KROCC	MSE	SROCC	PLCC	KROCC	MSE
MOS [34]	0.7537	0.7510	0.5638	0.0314	0.7848	0.7956	0.5889	0.0160	0.6584	0.6753	0.4756	0.0588
Zhi Li [35]	0.7537	0.7510	0.5638	0.0314	0.7848	0.7956	0.5889	0.0160	0.6584	0.6753	0.4756	0.0588
w/o score matrix	0.7858	0.7598	0.5800	0.0281	0.8151	0.8274	0.6181	0.0248	0.6821	0.6904	0.4955	0.0434
ours	0.8501	0.8220	0.6593	0.0181	0.8329	0.8492	0.6387	0.0085	0.7374	0.7474	0.5436	0.0293

Table 2. Performance comparison when the noise rate is set to different values. We conduct this test on the LIVE database and report the SROCC, PLCC, KROCC, and MSE results between the preprocessed quality scores and the true subjective scores of LIVE.

Datasets	LIVE				TID				CID			
Methods\Metrics	SROCC	PLCC	KROCC	MSE	SROCC	PLCC	KROCC	MSE	SROCC	PLCC	KROCC	MSE
MOS rate=0.6	0.8437	0.8407	0.6886	0.0193	0.8565	0.8617	0.6991	0.0099	0.7675	0.7794	0.5988	0.0336
Ours rate=0.6	0.8567	0.8339	0.6682	0.0190	0.8575	0.8713	0.6673	0.0083	0.8346	0.8350	0.6413	0.0201
MOS rate=0.8	0.8231	0.8131	0.6451	0.0246	0.8133	0.8230	0.6308	0.0134	0.7004	0.7141	0.5149	0.0465
Ours rate=0.8	0.8637	0.8429	0.6710	0.0169	0.8415	0.8550	0.6475	0.0082	0.7578	0.7595	0.5599	0.0277

If the dataset contains annotation information, we randomly sample the labeled scores of each image and average them to obtain a MOS with subjective bias as the training set. The original MOS is used as the test set. If the data set does not contain annotation information, we use the original MOS of each image as the mean, and the variance takes the given variance of the data set (or sets it to 0.2) as a Gaussian distribution to obtain a biased MOS.

Evaluation Criteria. The Pearson Linear Correlation Coefficient (PLCC), Spearman Rank Order Correlation Coefficient (SROCC), Kendall rank-order correlation coefficient (KROCC), and MSE are used to measure performance, as in previous work.

4.2. Detailed Implementation

Since perceptual similarity models have recently demonstrated great capability. We employ LPIPS and ResNet-50 as the backbone. The images in the training set are resized to 320, and randomly cropped to 320. The images in the test set are fed directly into our model with no data augmentation. All of this is done with the assumption that preprocessed images have the same score as the original. We use ResNet50 as the proposed method’s backbone network based on the general configuration of network structure in IQA fields. The model’s hyperparameter settings include learning rate = 0.01, SGD as the optimizer, epoch = 10, and training batch size = 16. The model was trained using a single GeForce RTX 3090 GPU.

4.3. Performance assessment and Comparison

In this section, experiments within individual standard IQA databases are conducted to evaluate the effectiveness of PSP. We discuss how to use single subjective labels to achieve performance under the labels obtained from many annotators. As the question has never been explored, we select the typical subjective assessment models, namely ResNet-50 [38] and LPIPS [33], and then test the accuracy of the models when using single subjective labels. To simulate the scenario of single labels, we replace all the labels in the datasets (except for CID2013). In other words, the bias label rate is 100 percent. Single subjective labels are generated via Gaussian Sampling from the labels processed by several annotators. Real single labels are obtained on CID2013 datasets, and the labels processed by several annotators are replaced randomly with the labels processed by a single annotator.

The experimental results in table 1 demonstrate the correlation between the score obtained by a small amount of score (just need one score in 1) and the real Ground Truth. The closer the first three indexes (SROCC, PLCC, and KROCC) are to one, the closer they are to Ground Truth, and the smaller the last is (MSE), the smaller the margin of error. Note that the comparison here is a comparison of the accuracy of the pre-processing method, and the quality of the labeling directly affects the performance of the model training. Compared with other mainstream crowdsourced subjective scores processing methods, such as MOS and Zhi Li’s, the proposed method is superior to others in SROCC, PLCC, KROCC, and MSE. The results are shown in Table 1. The accuracy of the models is significantly lower when trained with single subjective

Table 3. The performance impact of the dataset after the PSP subjective quality score preprocessing method on the quality assessment task. We conduct this test on the LIVE database and report SROCC, PLCC, KROCC, and MSE results between quality rating performance and LIVE true subjective scores. The quality evaluation model uses NIMA. The subjective bias scale was set to 1.0 and the bias variance to 0.2. We found that the scoring after PSP subjective quality score preprocessing can endow the quality evaluation model with better performance.

	SROCC	PLCC	KROCC	MSE
MOS-GT	0.9317	0.8988	0.7790	0.0112
MOS	0.8704	0.8653	0.6877	0.0238
PSP	0.9202	0.8966	0.7582	0.0178

labels than when trained with the labels processed in the original dataset (standard MOS). It indicates that the bias from the single subjective labels exerts a great negative impact on the accuracy of the models. However, the accuracy remains at similar levels when trained by the subjective labels processed by several annotators and when trained by the single subjective labels. When trained by single subjective labels, the proposed model is more accurate than other subjective biased models. we have the following observations.

4.4. Ablation studies

4.4.1. Different Noise Rates

Because there may be different annotations of labels for a single image in the case of true labeling, this paper calls it the noise rate. In order to simulate the scene and verify the effectiveness of the proposed PSP-IQA method, experiments are carried out under the noise rate of 0.6 and 0.8 respectively. The results are shown in Table 2, the proposed method outperforms the MOS method at different noise rates (e. g. 0.6,0.8,1), which verifies the universality of PSP-IQA.

4.4.2. Different EMA Weights

To evaluate the performance of the proposed model, we performed ablation experiments with respect to the EMA parameters of the model. In order to test the performance changes under different EMA weights, interval ablation experiments ranging from 0.2 to 1.0 were set up. The results are shown in Table 4.

4.4.3. Different variance and T

We performed ablation experiments to σ of the Gaussian distribution $p(y|x) \sim \mathcal{N}(u(x) + f(s), \sigma^2)$ and to the T of the

Table 4. Performance comparison when the EMA Weights is set to different values. We conduct this test on the LIVE database and report the SROCC, PLCC, KROCC, and MSE results between the preprocessed quality scores and the true subjective scores of LIVE. We highlight the best results in bold

	SROCC	PLCC	KROCC	MSE
EMA=0.2	0.4752	0.4324	0.3205	0.1247
EMA=0.4	0.5802	0.5492	0.4005	0.0758
EMA=0.8	0.7800	0.7484	0.5808	0.0271
EMA=0.9	0.8501	0.8220	0.6593	0.0181
EMA=1.0	0.7537	0.7510	0.5638	0.0314

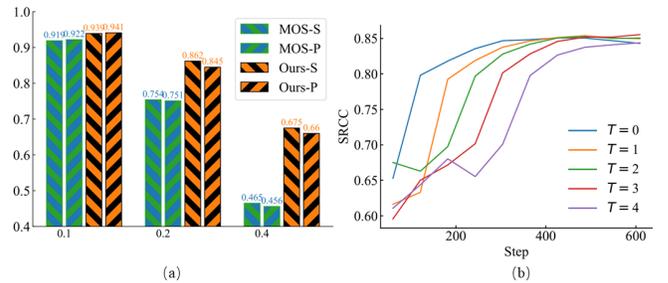


Fig. 3. (a) ResNet-50: Accuracies on LIVE with different variance σ .(b) ResNet-50: Accuracies on LIVE with different T .

EMA. The results are shown in Figure 3. As shown in figure 3 (a), the proposed PSP-IQA method is superior to the MOS method in different variance σ , and the larger the variance σ is, the more the performance gap between PSP-IQA and MOS method can be reflected, and the validity of the proposed method is verified. As shown in figure 3 (b), the convergence rate of PSP-IQA is different under different T , but the final convergence precision is close, which shows that the proposed method is insensitive to iterative parameters and robust.

4.5. Case Study

We conduct a case study to show the effectiveness of LPR-IQA. As show in Fig. 4, we found that calculating the perceptual similarity score between the two found a large gap with the actual annotated similarity score. In this way, we estimate a score matrix to fine-tune the score so that the score is more consistent with the actual distortion of the image. It is worth noting that our method fine-tunes the score by weighing the content of all images in the dataset and the distribution of scores, rather than the two images shown. The marked score of the image we selected is very different from the real score, and the corresponding reference image searched by NSS is



Fig. 4. Given a biased MOS (below the image), our method searches for images with similar perceptual similarities via NSS and gets the finetune score.

also the same. There is a gap between our perceptual similarity score (0.2905) and the annotated similarity score (0.5441). However, the perceptual similarity score is correlated with the whole dataset with higher confidence. Based on the perceptual similarity score, we were able to fine-tune the image score.

5. CONCLUSIONS AND FUTURE WORK

This paper proposes a new IQA-perceptual similarity processing method, which is a simple and effective image-perceptual similarity-based IQA data preprocessing method PSP-IQA, suitable for IQA tasks. PSP-IQA alleviates the subjective bias problem when there are few annotators. This method can consider not only the existing annotation labels, but also the perceptual similarity relationship of images in the data set. In our approach, we assume that the subjective annotation process is a scoring process based on latent reference images, relative to the perceived similarity of another image. We propose a conditional probability model based on LPISP and nearest neighbor search (NSS) to model the subjective scoring process. And fine-tune the existing labels based on EMA. On the real dataset CID2013 and the popular IQA datasets LIVE and TID, it is verified that our method can effectively alleviate the bias of subjective scoring.

6. REFERENCES

[1] Kede Ma, Wentao Liu, Kai Zhang, Zhengfang Duanmu, Zhou Wang, and Wangmeng Zuo, “End-to-end blind image quality assessment using deep neural networks,” *IEEE Transactions on Image Processing*, vol. 27, no. 3, pp. 1202–1213, 2017.

[2] Sebastian Bosse, Dominique Maniry, Klaus-Robert

Müller, Thomas Wiegand, and Wojciech Samek, “Deep neural networks for no-reference and full-reference image quality assessment,” *IEEE Transactions on image processing*, vol. 27, no. 1, pp. 206–219, 2017.

[3] Weixia Zhang, Kede Ma, Jia Yan, Dexiang Deng, and Zhou Wang, “Blind image quality assessment using a deep bilinear convolutional neural network,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 1, pp. 36–47, 2018.

[4] Hanhe Lin, Vlad Hosu, and Dietmar Saupe, “Kadid-10k: A large-scale artificially distorted iqa database,” in *2019 Eleventh International Conference on Quality of Multimedia Experience (QoMEX)*. IEEE, 2019, pp. 1–3.

[5] Debarati Kundu, Deepti Ghadiyaram, Alan C Bovik, and Brian L Evans, “Large-scale crowdsourced study for tone-mapped hdr pictures,” *IEEE Transactions on Image Processing*, vol. 26, no. 10, pp. 4725–4740, 2017.

[6] Toni Virtanen, Mikko Nuutinen, Mikko Vaahteranoksa, Pirkko Oittinen, and Jukka Häkkinen, “Cid2013: A database for evaluating no-reference image quality assessment algorithms,” *IEEE Transactions on Image Processing*, vol. 24, no. 1, pp. 390–402, 2014.

[7] Andela Zaric, Nenad Tatalovic, Nikolina Brajkovic, Hrvoje Hlevnjak, Matej Loncaric, Emil Dumic, and Sonja Grgic, “Vcl@ fer image quality assessment database,” in *Proceedings ELMAR-2011*. IEEE, 2011, pp. 105–110.

[8] Telephone Installations and Local Line, “Itu-tp. 910,” *Subjective video quality assessment methods for multimedia applications, Recommendation ITU-T*, p. 910.

[9] IT Union, “Methods for the subjective assessment of video quality audio quality and audiovisual quality of internet video and distribution quality television in any environment,” *SERIES P: TERMINALS AND SUBJECTIVE AND OBJECTIVE ASSESSMENT METHODS*, 2016.

[10] I BT, “Methodologies for the subjective assessment of the quality of television images, document recommendation itu-r bt. 500-14 (10/2019),” *ITU, Geneva, Switzerland*, 2020.

[11] Jupo Ma, Jinjian Wu, Leida Li, Weisheng Dong, Xue-mei Xie, Guangming Shi, and Weisi Lin, “Blind image quality assessment with active inference,” *IEEE Transactions on Image Processing*, vol. 30, pp. 3650–3663, 2021.

[12] Lei Wang, Qingbo Wu, King Ngai Ngan, Hongliang Li, Fanman Meng, and Linfeng Xu, “Blind tone-mapped image quality assessment and enhancement via disentangled representation learning,” in *2020 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*. IEEE, 2020, pp. 1096–1102.

[13] Junyong You and Jari Korhonen, “Transformer for image quality assessment,” in *2021 IEEE International*

- Conference on Image Processing (ICIP)*. IEEE, 2021, pp. 1389–1393.
- [14] Huiyu Duan, Guangtao Zhai, Xiongkuo Min, Yucheng Zhu, Yi Fang, and Xiaokang Yang, “Perceptual quality assessment of omnidirectional images,” in *2018 IEEE international symposium on circuits and systems (IS-CAS)*. IEEE, 2018, pp. 1–5.
- [15] Wei Sun, Xiongkuo Min, Guangtao Zhai, Ke Gu, Huiyu Duan, and Siwei Ma, “Mc360iqa: A multi-channel cnn for blind 360-degree image quality assessment,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, no. 1, pp. 64–77, 2019.
- [16] Vamsi Kiran Adhikarla, Marek Vinkler, Denis Sumin, Rafał Mantiuk, Karol Myszkowski, Hans-Peter Seidel, and Piotr Didyk, “Towards a quality metric for dense light fields,” in *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [17] Xiongkuo Min, Guangtao Zhai, Ke Gu, Xiaokang Yang, and Xinpeng Guan, “Objective quality evaluation of dehazed images,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 20, no. 8, pp. 2879–2892, 2018.
- [18] Xiongkuo Min, Guangtao Zhai, Ke Gu, Yucheng Zhu, Jiantao Zhou, Guodong Guo, Xiaokang Yang, Xinpeng Guan, and Wenjun Zhang, “Quality evaluation of image dehazing methods using synthetic hazy images,” *IEEE Transactions on Multimedia*, vol. 21, no. 9, pp. 2319–2333, 2019.
- [19] Devansh Arpit, Stanislaw Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al., “A closer look at memorization in deep networks,” in *International conference on machine learning*. PMLR, 2017, pp. 233–242.
- [20] Aritra Ghosh, Himanshu Kumar, and PS Sastry, “Robust loss functions under label noise for deep neural networks,” in *Proceedings of the AAAI conference on artificial intelligence*, 2017, vol. 31.
- [21] Daiki Tanaka, Daiki Ikami, Toshihiko Yamasaki, and Kiyoharu Aizawa, “Joint optimization framework for learning with noisy labels,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 5552–5560.
- [22] Jacob Goldberger and Ehud Ben-Reuven, “Training deep neural-networks using a noise adaptation layer,” 2016.
- [23] Xiaobo Xia, Tongliang Liu, Bo Han, Mingming Gong, Jun Yu, Gang Niu, and Masashi Sugiyama, “Sample selection with uncertainty of losses for learning with noisy labels,” *arXiv preprint arXiv:2106.00445*, 2021.
- [24] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals, “Understanding deep learning (still) requires rethinking generalization,” *Communications of the ACM*, vol. 64, no. 3, pp. 107–115, 2021.
- [25] Zhilu Zhang and Mert Sabuncu, “Generalized cross entropy loss for training deep neural networks with noisy labels,” *Advances in neural information processing systems*, vol. 31, 2018.
- [26] Xiong Zhou, Xianming Liu, Chenyang Wang, Deming Zhai, Junjun Jiang, and Xiangyang Ji, “Learning with noisy labels via sparse regularization,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 72–81.
- [27] Sangchul Hahn and Heeyoul Choi, “Self-knowledge distillation in natural language processing,” in *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2019)*, 2019, pp. 423–430.
- [28] Akhilesh Gotmare, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher, “A closer look at deep learning heuristics: Learning rate restarts, warmup and distillation,” in *International Conference on Learning Representations*, 2019.
- [29] Yisen Wang, Xingjun Ma, Zaiyi Chen, Yuan Luo, Jinfeng Yi, and James Bailey, “Symmetric cross entropy for robust learning with noisy labels,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 322–330.
- [30] Honglei Zhuang and Joel Young, “Leveraging in-batch annotation bias for crowdsourced active learning,” in *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining*, 2015, pp. 243–252.
- [31] Christoph Hube, Besnik Fetahu, and Ujwal Gadiraju, “Understanding and mitigating worker biases in the crowdsourced collection of subjective judgments,” in *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, 2019, pp. 1–12.
- [32] Guoliang Li, Jiannan Wang, Yudian Zheng, and Michael J Franklin, “Crowdsourced data management: A survey,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 28, no. 9, pp. 2296–2319, 2016.
- [33] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang, “The unreasonable effectiveness of deep features as a perceptual metric,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.
- [34] Robert C Streijl, Stefan Winkler, and David S Hands, “Mean opinion score (mos) revisited: methods and applications, limitations and alternatives,” *Multimedia Systems*, vol. 22, no. 2, pp. 213–227, 2016.
- [35] Zhi Li, Christos G Bampis, Lucjan Janowski, and Ioannis Katsavounidis, “A simple model for subject behavior in subjective experiments,” *Electronic Imaging*, vol. 2020, no. 11, pp. 131–1, 2020.
- [36] Hamid R Sheikh, Muhammad F Sabir, and Alan C Bovik, “A statistical evaluation of recent full reference image quality assessment algorithms,” *IEEE Transac-*

tions on image processing, vol. 15, no. 11, pp. 3440–3451, 2006.

- [37] Nikolay Ponomarenko, Lina Jin, Oleg Ieremeiev, Vladimir Lukin, Karen Egiazarian, Jaakko Astola, Benoit Vozel, Kacem Chehdi, Marco Carli, Federica Battisti, et al., “Image database tid2013: Peculiarities, results and perspectives,” *Signal processing: Image communication*, vol. 30, pp. 57–77, 2015.
- [38] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.