

# Opinion-Unaware Blind Image Quality Assessment using Multi-Scale Deep Feature Statistics

Zhangkai Ni, *Member, IEEE*, Yue Liu, Keyan Ding, Wenhan Yang, *Member, IEEE*, Hanli Wang, *Senior Member, IEEE*, Shiqi Wang, *Senior Member, IEEE*

**Abstract**—Deep learning-based methods have significantly influenced the blind image quality assessment (BIQA) field, however, these methods often require training using large amounts of human rating data. In contrast, traditional knowledge-based methods are cost-effective for training but face challenges in effectively extracting features aligned with human visual perception. To bridge these gaps, we propose integrating deep features from pre-trained visual models with a statistical analysis model into a Multi-scale Deep Feature Statistics (MDFS) model for achieving opinion-unaware BIQA (OU-BIQA), thereby eliminating the reliance on human rating data and significantly improving training efficiency. Specifically, we extract patch-wise multi-scale features from pre-trained vision models, which are subsequently fitted into a multivariate Gaussian (MVG) model. The final quality score is determined by quantifying the distance between the MVG model derived from the test image and the benchmark MVG model derived from the high-quality image set. A comprehensive series of experiments conducted on various datasets show that our proposed model exhibits superior consistency with human visual perception compared to state-of-the-art BIQA models. Furthermore, it shows improved generalizability across diverse target-specific BIQA tasks. Our code is available at: <https://github.com/eezkn1/MDFS>

**Index Terms**—Blind image quality assessment, multivariate Gaussian fitting, multi-scale deep features, feature statistics

## I. INTRODUCTION

IMAGE quality assessment (IQA) is a critical and fundamental research topic in the field of computer vision due to its extensive applicability in various tasks, including image compression [1], image super-resolution [2], and image enhancement [3], [4]. Since the primary recipients of images are humans, the image quality scores derived from subjective quality assessment experiments are highly reliable,

however, conducting such experiments is expensive and time-consuming. Therefore, numerous IQA models have been proposed in the past half-century to predict image quality that is highly consistent with the human visual system (HVS).

Existing IQA models can be divided into three categories according to the use of reference images: full-reference IQA (FR IQA), reduced-reference IQA (RR IQA), and blind IQA (BIQA). Unlike reference-based methods, BIQA alleviates the need for direct comparisons against pristine reference images, which is often not feasible in real-world scenarios. Human observers are capable of assessing image quality even in the absence of reference images, implying that the HVS is good at perceiving perceptual characteristics closely associated with natural image quality. Our research scope is dedicated to BIQA models, with the primary objective of developing an effective and reliable method for predicting the quality score of distorted images in the absence of corresponding reference images.

BIQA methods can be divided into opinion-aware BIQA (OA-BIQA) and opinion-unaware BIQA (OU-BIQA), depending on whether relying on subjective scores for training. Currently, most research efforts are focused on OA-BIQA models, which require training with datasets containing human-rated quality labels. However, applying deep learning to IQA is challenging due to the typically small sizes of IQA datasets [5], [6], which increases the risk of overfitting in deep learning-based IQA models [7], [8]. OU-BIQA models evaluate image quality solely based on the visual features and characteristics [9], [10]. The advantage of these methods is that no subjective scoring is required during the training process, thereby reducing potential gaps in the evaluation process caused by the subjectivity of different datasets, which leads to improved generalization and robustness. Consequently, our work emphasizes developing a robust OU-BIQA model. Traditional OU-BIQA methods [11], [12] have the advantage of low training costs but leave room for further improvement in performance. This is primarily because these approaches may not fully capture the intrinsic image characteristics aligned with human perception when assessing image quality. Therefore, investigating the extraction of image features aligned with HVS perception is critical for the OU-BIQA algorithm. Deep networks enable the model to learn intrinsic representations and automatically capture important subtle image characteristics from the data [13]–[15]. A straightforward alternative approach is to utilize a pre-trained deep network as a feature extractor, initially trained on a large-scale dataset not specifically related to the target data (*i.e.*, IQA images), thus enabling the network to acquire abstract and more universally applicable features.

This work was supported in part by the National Natural Science Foundation of China under Grant 62201387, Grant 62371343, and Grant 62301480, in part by the Shanghai Pujiang Program under Grant 22PJ1413300, and in part by the Fundamental Research Funds for the Central Universities. (*Corresponding authors: Hanli Wang and Keyan Ding*)

Zhangkai Ni and Hanli Wang are with the Department of Computer Science and Technology, Key Laboratory of Embedded System and Service Computing (Ministry of Education), and Shanghai Institute of Intelligent Science and Technology, Tongji University, Shanghai 200092, China (e-mail: zkni@tongji.edu.cn; hanliwang@tongji.edu.cn).

Yue Liu and Shiqi Wang are with the Department of Computer Science, City University of Hong Kong, Hong Kong 999077 (e-mail: yliu724-c@my.cityu.edu.hk; shiqiwan@cityu.edu.hk).

Keyan Ding is with ZJU-Hangzhou Global Scientific and Technological Innovation Center, Zhejiang University, Hangzhou, Zhejiang 311200, China (e-mail: dingkeyan@zju.edu.cn).

Wenhan Yang is with PengCheng Laboratory, Shenzhen, Guangdong 518066, China. (e-mail: yangwh@pcl.ac.cn).

Considering the complementary advantages of these two kinds of methods, we aim to integrate the robustness of the traditional statistical analysis model and the universally applicable features provided by a pre-trained deep model. Thus, our method aims to leverage the power of DNNs to extract multi-scale features from images, and then employ statistical analysis algorithms to further analyze and aggregate the extracted features, which are integrated into the proposed Multi-scale Deep Feature Statistics (MDFS) model for OU-BIQA. Specifically, we first utilize a pre-trained network trained on a large-scale dataset agnostic to IQA images to generate multi-scale feature maps, and subsequently down-sample and concatenate these feature maps to build feature maps with rich context and semantics. In the statistical data analysis stage, we calculate the mean and variance of these feature maps and utilize a multivariate Gaussian (MVG) model to model its distribution. The final MDFS index is computed by quantifying the similarity between the MVG model derived from the features of the testing image and the benchmark MVG model derived from features of a training image dataset containing only high-quality images. The main contributions of our works are summarized as follows:

- We proposed the Multi-scale Deep Feature Statistic (MDFS) model for OU-BIQA, which integrates deep features extracted from DNNs into a statistical data analysis model and derives their feature distribution using an MVG model. To our knowledge, this is the first attempt to combine deep features with data distribution of traditional methods in the context of the OU-BIQA model.
- We designed a statistical data analysis model to extract discriminative information from deep feature maps and integrate this information into the MVG model. Our work serves as a bridge to seamlessly incorporate deep features into the established traditional statistical modeling, thus paving a new way for the systematic fusion of the methods of these two categories.
- The proposed MDFS outperforms state-of-the-art OU-BIQA methods in terms of cost-effectiveness during training and superior performance across diverse datasets. Moreover, it can be easily generalized to various target-specific BIQA tasks.

The rest of this paper is organized as follows. Section II outlines the related work. Section III introduces the proposed MDFS model in detail. Section IV presents a comprehensive comparison and analysis of experimental results. Finally, Section V draws the conclusion.

## II. RELATED WORKS

This section provides an overview of related BIQA methods, highlighting their scalability and flexibility advantages over various FR IQA methods specialized for different image types, such as stereoscopic images [16], [17], retargeted images [18], and multi-exposure fused images [19].

### A. Traditional BIQA Methods

Early BIQA methods are designed to assess image degradation caused by specific distortion types. However, given

the typically unknown and diverse distortion types present in a single image, there is a demand for generalized BIQA approaches capable of adapting to various complex scenarios. Natural scene statistics (NSS) based BIQA methods leverage the assumption that NSS is closely related to image degradation [20], extracting diverse NSS features in spatial or transformed domains to estimate distribution disparity between reference and distorted images.

In the spatial domain, Mittal *et al.* [11] developed the natural image quality evaluator (NIQE), where a set of image patches are selected and the corresponding statistics are regarded as quality-aware features. The multivariate Gaussian (MVG) model is used to estimate the global distribution of natural images, which is compared with the corresponding MVG model of the distorted image to predict the final quality score. Zhang *et al.* [21] enhanced the NIQE by incorporating color, gradient, and frequency characteristics in the feature extraction stage, resulting in the Integrated Local NIQE (IL-NIQE). To reduce the feature dimension in ILNIQE, Liu *et al.* [22] implies sparse representation in the proposed structure, naturalness, and perception quality-driven NIQE (SNP-NIQE). They further proposed the natural scene statistics and perceptual characteristics-based quality index (NPQI) [23] by applying the local binary pattern map and the locally mean subtracted and contrast normalized (MSCN) coefficients of the image to extract the NSS features. Furthermore, Xue *et al.* [24] proposed a quality-aware clustering (QAC) method by learning a set of quality centroids, which is used to estimate the quality of each patch of the input image. Later, Venkatanath *et al.* [12] introduced a visual attention strategy and developed a perception-based image quality evaluator (PIQE) to improve the assessment accuracy. Wu *et al.* [25] designed a local pattern statistics index (LPSI) by modifying the statistics of the feature extracted from the local binary pattern.

In addition to directly extracting NSS features in the spatial domain, BIQA methods have explored frequency domain features, such as those in the discrete cosine transform (DCT) domain and wavelet domain. Moorthy *et al.* [26] modeled the wavelet coefficients and further identified the distortion type of the image using a support vector machine (SVM) to produce the final quality score. Saad *et al.* [27] applied a generalized Gaussian distribution (GGD) model to predict image quality based on features extracted from the DCT coefficients. Furthermore, Wu *et al.* [28] combined features from multiple domains and color channels in the proposed type classification and label transfer (TCLT) model to estimate perceptual image quality.

### B. Learning-based BIQA Methods

In recent years, convolutional neural networks (CNNs) have made significant strides in the field of BIQA, with techniques including graph convolutional neural networks [29], continual learning approaches [30], and attention mechanisms [31], [32]. These methods can be classified into supervised (opinion-aware) and unsupervised (opinion-unaware) BIQA algorithms, based on the availability of subjective scores for training. For opinion-aware OA-BIQA methods, Kim *et al.* [33] applied an

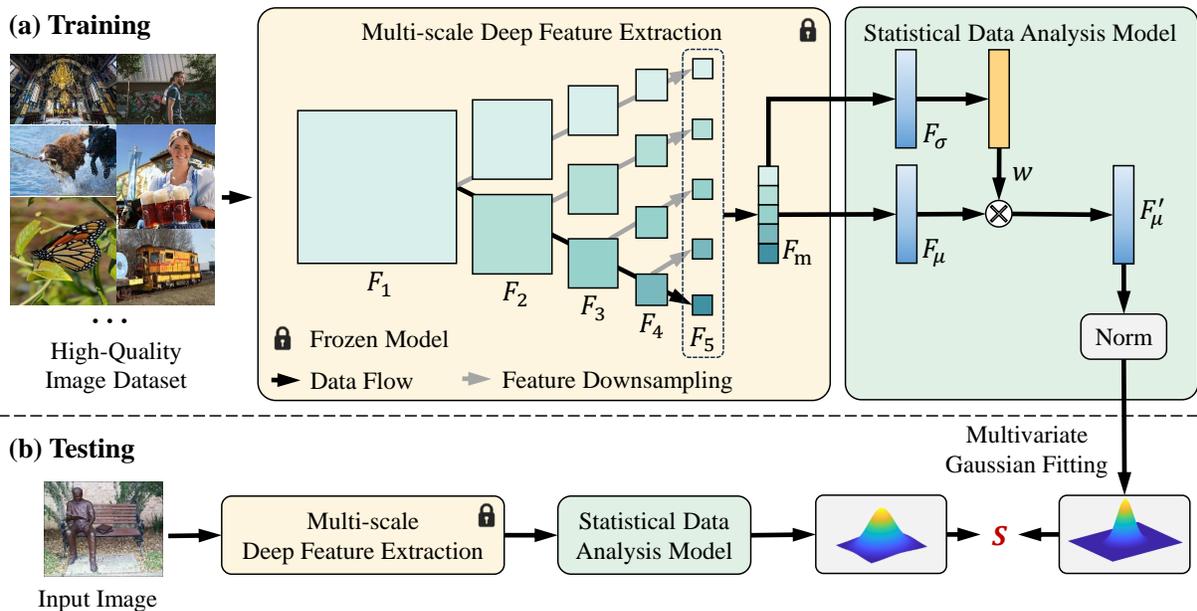


Fig. 1: Overview of the proposed MDFS model: (a) Training phase: This process involves fitting a benchmark multivariate Gaussian (MVG) model from a set of high-quality images, including a frozen multi-scale deep feature extraction module (e.g., ResNet, VGG, and EfficientNet), a statistical data analysis model, and an MVG fitting model. (b) Testing phase: The process of assessing the quality of a test image involves calculating the final quality score by measuring the distance between an MVG model fitted using the test image features and the benchmark MVG model obtained in the training phase.

FR IQA method as an auxiliary task to guide the network to learn quality map, where intermediate features are utilized to predict the final quality score. Guan *et al.* [8] proposed to extract features of each patch and then use a local regression module to generate local responses and weights, which are combined for final quality prediction. Ma *et al.* [7] built an end-to-end network where the first sub-network classifies distortion types to assist the second sub-network in severity assessment. To combine more IQA datasets for BIQA model training, Zhang *et al.* [30] introduced a new head for new datasets in the training process, along with existing heads for predicting image quality, where k-means clustering is used to generate the weighting map. Besides, graph representation has been applied in BIQA [29], where the node and edge of the learned graph represent the distortion types and distortion levels. Recently, Su *et al.* [31] took the attention mechanism into the BIQA model for better content understanding. Yang *et al.* [34] also introduced a transformer in the BIQA task.

OA-BIQA methods can achieve state-of-the-art performance on specific datasets with the same distribution as the training dataset, but their generalizability to new datasets has been discredited due to potential overfitting issues [30]. In contrast, OU-BIQA algorithms make it easier to establish the training datasets and can adapt to adjustments based on new, unlabeled datasets, making them more useful for various applications. As a result, various researchers have shifted their focus toward OU-BIQA algorithms. Ma *et al.* [13] proposed to generate a set of distorted image pairs, where higher quality images are identified by multiple IQA methods. This generated dataset is later used for image quality estimation. In their subsequent

work [35], they further improved this model by combining the training data into a quality discriminable image pair (DIP) format, which is subsequently fed into a pairwise learning-to-rank algorithm for quality measurement. Recently, Chen *et al.* [36] proposed a self-supervised strategy, where the quality-aware information is learned from a patch prediction framework based on contrastive learning. Babu *et al.* [37] proposed a self-supervised method, where mutual information bounds are used to separate content information from image patches, focusing on content-independent image quality.

### III. MULTI-SCALE DEEP FEATURE STATISTIC MODEL

#### A. Overview

Previous OA-BIQA methods suffer from the risk of overfitting due to the limited data while conventional OU-BIQA fails to obtain satisfactory performance caused by the lack of intrinsic image features. Our goal is to build a robust and efficient OU-IQA method by applying statistical analysis models for universally applicable features extracted from a pre-trained deep network to simultaneously inherit high performance in terms of robustness and training efficiency. We innovatively leverage multi-scale features learned by a deep neural network to fit a multivariate Gaussian (MVG) model. The framework of our proposed MDFS for OU-BIQA is shown in Figure 1, which consists of training and testing phases, outlined as follows:

- In the training phase, image features are initially extracted from a pre-trained neural network using **Multi-Scale Deep Feature Extraction**, followed by analysis with a

**Statistical Data Analysis Model**, and then fitted to the benchmark **Multivariate Gaussian** (MVG) model.

- In the testing phase, the MVG model of the distorted image is obtained through a process similar to the training phase. The final quality score of the test image is determined by evaluating the distance between the MVG model of the test image and the benchmark MVG model.

Each component of the two phases will be detailed in the following subsections, respectively.

### B. Multi-scale Deep Feature Extraction

In recent years, deep learning has significantly improved the accuracy and efficiency of numerous computer vision tasks, with various classic network architectures proposed, such as ViT [38], VGG [39], ConvNet [40], ResNet [41], and EfficientNet [42], serving as the foundation for various downstream tasks. The multi-scale deep feature extraction modal extracts multi-scale features in a pyramid form for better feature representation. Specifically, given an input image  $I_{in}$ , the pre-trained vision modal (*e.g.*, VGG, ResNet, and EfficientNet) is first used to extract image features. The outputs of the first layer to the fifth layer of the model are specially extracted and represented as  $F_i (i = 1, 2, 3, 4, 5)$ . Features of the first three layers are downsampled to unify multiple features of different scales, which is defined as,

$$F_m = F_5 \circ D(F_4 \circ D(F_3 \circ D(F_2 \circ D(F_1))))), \quad (1)$$

where  $D(\cdot)$  indicates the downsampling operation and  $\circ$  is concatenation. The downsampling is implemented using the same non-learnable convolutional layer with a stride of 2 to halve the scale of feature maps, and reflection padding is applied to mitigate edge artifacts. Note that the multi-scale deep feature extraction is not learnable and the kernel used in the downsampling layer remains fixed during the training and testing stage.

### C. Statistical Data Analysis Model

In traditional BIQA methods [11], [23], NSS features are widely explored to identify a perceptual feature extraction function that closely approximates the sensitivity of the HVS. However, how to perform statistical analysis for learning-based features has received little attention. In this section, a statistical data analysis model is proposed for further data distillation and analysis based on extracted deep features. In this model, the non-learnable convolution operation is utilized to estimate the local mean value of the feature map, which is subsequently screened before normalization. Specifically, given an input deep feature  $F_m \in \mathbb{R}^{C \times H \times W}$ , where  $C$ ,  $H$ , and  $W$  denotes the number of channels, height, and width respectively, a Gaussian filter with dynamic window size  $s_w$  is applied on the feature to estimate the local mean value of the feature map. Therefore, statistical features, including the mean and standard deviation of  $F_m$ , are generated as follows,

$$F_\mu = \text{conv}(F_m, s_w), \quad (2)$$

$$F_\sigma = \text{mean}_c \sqrt{\text{conv}(F_m^2, s_w)}, \quad (3)$$

where  $\text{conv}(\cdot)$  and  $\text{mean}_c(\cdot)$  denote the convolutional operation with the stride of 1 and the averaging operation in the channel dimension, respectively. To enable the filter to adapt to feature maps of different sizes, a dynamic window size calculation method is employed as follows,

$$s_w = \max(3, 1 + 2 \cdot (\min(H, W) // 2^k)), \quad (4)$$

where  $//$  represents the remainder operation and  $k$  is empirically set to 5 in this work. Since the feature map is obtained by concatenating features from five different scales, we normalize  $F_\mu$  along the channel dimension of each layer to obtain the final features for:

$$F'_\mu = \text{norm}_c(F_\mu), \quad (5)$$

where  $\text{norm}_c(\cdot)$  represents the normalization operation along the channel dimension.

Considering that the HVS is more sensitive to information with higher contrast [43], we utilize standard deviation as a measure of contrast and highlight local regions exhibiting greater standard deviation. The weighting map is defined as:

$$w = \frac{1}{1 + e^{-(F_\sigma - F_\sigma^\mu)/(F_\sigma^\sigma + \delta)}}, \quad (6)$$

where  $F_\sigma^\mu$  and  $F_\sigma^\sigma$  are the mean and standard deviation of the  $F_\sigma$ , respectively.  $\delta$  is a small positive number (*i.e.*,  $\delta = 1 \times e^{-12}$ ) to avoid the denominator being equal to 0. The weighting map of the distorted image is adopted in the testing stage to calculate the weighted quality score.

### D. Multivariate Gaussian Model

The MVG model has been extensively utilized to model the joint probability distribution of a vector of random variables, with each variable following a normal distribution. Herein, we use the MVG model to estimate the joint probability distribution of a set of training images. Let  $\vec{X} = [x_1, \dots, x_n]$ , ( $n = 1, 2, 3, \dots$ ) represent the statistical features obtained from  $n$  high-quality images. Assuming that these features represent independent samples from an  $l$ -dimensional MVG distribution, the MVG model learned through Maximum Likelihood Estimation (MLE) can be expressed as follows:

$$p(\vec{X}) = \frac{1}{\sqrt{(2\pi)^l |\Sigma|}} \cdot e^{-\frac{1}{2}(\vec{X} - \vec{\mu})^T \Sigma^{-1}(\vec{X} - \vec{\mu})}, \quad (7)$$

where  $l$  is the dimension of the learned statistical features, while  $\vec{\mu}$  and  $\Sigma$  denote the mean and covariance matrix of the estimated MVG model, respectively.

### E. Quality Calculation

The proposed quality index is computed by measuring the distance between the MVG model fitted using the features of the testing image and the benchmark MVG model fitted using the features from a high-quality image set. All the deep features are utilized in the statistical data analysis model for the testing image without removing low contrast features as suggested in [11]. Specifically, the quality score of the test image is calculated by the Mahalanobis distance between two MVGs as suggested in [11], [21], [22]. The covariance matrix

TABLE I: Performance comparisons of different OU-BIQA models on ten public datasets. The top three performers are marked in bold red, blue, and black, respectively.

Metrics	Datasets	NIQE	QAC	PIQE	LPSI	ILNIQE	dipIQ	SNP-NIQE	NPQI	ContentSep	MDFS (Ours)
SROCC	LIVE	0.9062	0.8683	0.8398	0.8181	0.8975	<b>0.9378</b>	0.9073	<b>0.9108</b>	0.7478	<b>0.9361</b>
	CSIQ	0.6191	0.4804	0.5120	0.5218	<b>0.8045</b>	0.5191	0.6090	<b>0.6341</b>	0.5871	<b>0.7774</b>
	TID2013	0.3106	0.3719	0.3636	0.3949	<b>0.4938</b>	<b>0.4377</b>	0.3329	0.2804	0.2530	<b>0.5363</b>
	KADID	0.3779	0.2394	0.2372	0.1478	<b>0.5406</b>	0.2977	0.3719	0.3909	<b>0.5060</b>	<b>0.5983</b>
	MDLIVE	0.7728	0.4116	0.3862	0.2717	<b>0.8778</b>	0.6678	<b>0.7822</b>	<b>0.8100</b>	0.4285	0.7579
	MDIVL	0.5656	0.5524	0.5319	0.5736	0.6237	<b>0.7131</b>	<b>0.6252</b>	0.6139	0.2582	<b>0.7890</b>
	KonIQ	0.5300	0.3397	0.2452	0.2239	0.5057	0.2375	<b>0.6284</b>	0.6132	<b>0.6401</b>	<b>0.7333</b>
	CLIVE	0.4495	0.2258	0.2325	0.0832	0.4393	0.2089	0.4654	<b>0.4752</b>	<b>0.5060</b>	<b>0.4821</b>
	CID2013	0.6589	0.0299	0.0448	0.3229	0.3062	0.3776	<b>0.7159</b>	<b>0.7698</b>	0.6116	<b>0.8571</b>
	SPAQ	0.3105	0.4397	0.2317	0.0001	<b>0.6959</b>	0.2189	0.5402	0.5999	<b>0.7084</b>	<b>0.7408</b>
	$AVG_D$	0.5501	0.3959	0.3625	0.3358	<b>0.6185</b>	0.4616	0.5978	<b>0.6098</b>	0.5247	<b>0.7208</b>
	$AVG_W$	0.4226	0.3562	0.2706	0.1760	<b>0.5851</b>	0.2987	0.5164	0.5322	<b>0.5815</b>	<b>0.6854</b>
KROCC	LIVE	0.7275	0.6738	0.6367	0.6175	0.7123	<b>0.7806</b>	0.7353	<b>0.7421</b>	0.5456	<b>0.7709</b>
	CSIQ	0.4520	0.3452	0.3687	0.3736	<b>0.6109</b>	0.3963	0.4492	<b>0.4819</b>	0.4057	<b>0.5823</b>
	TID2013	0.2114	0.2575	0.2554	0.2734	<b>0.3491</b>	<b>0.3016</b>	0.2285	0.1960	0.1676	<b>0.3824</b>
	KADID	0.2624	0.1660	0.1657	0.1004	<b>0.3808</b>	0.2133	0.2584	0.2732	<b>0.3426</b>	<b>0.4238</b>
	MDLIVE	0.5799	0.2903	0.2715	0.2071	<b>0.6880</b>	0.4872	<b>0.5923</b>	<b>0.6196</b>	0.2917	0.5623
	MDIVL	0.3934	0.3751	0.3642	0.3998	0.4383	<b>0.5034</b>	<b>0.4444</b>	0.4359	0.1641	<b>0.5911</b>
	KonIQ	0.3679	0.2302	0.1649	0.1504	0.3504	0.1594	<b>0.4434</b>	0.4310	<b>0.4529</b>	<b>0.5344</b>
	CLIVE	0.3064	0.1514	0.1561	0.0523	0.2984	0.1395	0.3162	<b>0.3256</b>	<b>0.3450</b>	<b>0.3274</b>
	CID2013	0.4675	0.0178	0.0394	0.2168	0.2100	0.2614	<b>0.5156</b>	<b>0.5655</b>	0.4378	<b>0.6706</b>
	SPAQ	0.2059	0.3001	0.1560	0.0006	<b>0.4930</b>	0.1454	0.3686	0.4137	<b>0.5069</b>	<b>0.5347</b>
	$AVG_D$	0.3974	0.2807	0.2579	0.2392	<b>0.4531</b>	0.3388	0.4352	<b>0.4485</b>	0.3660	<b>0.5380</b>
	$AVG_W$	0.2932	0.2456	0.1869	0.1216	<b>0.4145</b>	0.2093	0.3619	0.3749	<b>0.4075</b>	<b>0.4966</b>
PLCC	LIVE	0.9041	0.8625	0.8197	0.7859	0.9022	<b>0.9295</b>	<b>0.9060</b>	<b>0.9161</b>	0.4639	0.8558
	CSIQ	0.6901	0.5934	0.6279	0.6950	<b>0.7232</b>	<b>0.7009</b>	0.6962	0.6479	0.3632	<b>0.7907</b>
	TID2013	0.3789	0.4190	0.4615	0.4594	<b>0.5090</b>	<b>0.4746</b>	0.4055	0.4000	0.2203	<b>0.6242</b>
	KADID	0.3883	0.3088	0.2887	0.3348	<b>0.5341</b>	0.3832	<b>0.4212</b>	0.3401	0.3568	<b>0.5939</b>
	MDLIVE	0.8378	0.4149	0.3778	0.3727	<b>0.8923</b>	0.7241	<b>0.8525</b>	<b>0.8454</b>	0.3524	0.8226
	MDIVL	0.5650	0.5713	0.5142	0.5715	0.5697	<b>0.7252</b>	<b>0.6393</b>	0.6013	0.2311	<b>0.7953</b>
	KonIQ	0.4835	0.2906	0.2061	0.1064	0.4963	0.3773	<b>0.6222</b>	0.6139	<b>0.6274</b>	<b>0.7123</b>
	CLIVE	0.4939	0.2841	0.3144	0.2521	0.5033	0.3163	<b>0.5199</b>	0.4920	<b>0.5130</b>	<b>0.5364</b>
	CID2013	0.6712	0.0981	0.1072	0.4439	0.4267	0.3829	<b>0.7260</b>	<b>0.7772</b>	0.6368	<b>0.8717</b>
	SPAQ	0.2639	0.4497	0.2488	0.1183	<b>0.6371</b>	0.2239	0.5469	0.6155	<b>0.6648</b>	<b>0.7177</b>
	$AVG_D$	0.5677	0.4292	0.3966	0.4140	0.6194	0.5238	<b>0.6336</b>	<b>0.6249</b>	0.4430	<b>0.7321</b>
	$AVG_W$	0.4090	0.3735	0.2914	0.2441	<b>0.5661</b>	0.3697	<b>0.5399</b>	0.5340	0.5127	<b>0.6803</b>
RMSE	LIVE	11.6733	13.8258	15.6508	16.8932	11.7834	<b>10.0761</b>	<b>11.5643</b>	<b>10.9529</b>	24.2043	14.1344
	CSIQ	0.1900	0.2113	0.2043	0.1888	<b>0.1813</b>	<b>0.1873</b>	0.1885	0.2000	0.2631	<b>0.1607</b>
	TID2013	1.1472	1.1256	1.0998	1.1011	<b>1.0670</b>	<b>1.0911</b>	1.1332	1.1362	1.2092	<b>0.9685</b>
	KADID	0.9977	1.0297	1.0365	1.0201	<b>0.9153</b>	1.0000	<b>0.9819</b>	1.0181	1.0114	<b>0.8710</b>
	MDLIVE	10.3244	17.2073	17.5107	17.5497	<b>8.5379</b>	13.0432	<b>9.8857</b>	<b>10.1012</b>	17.6986	10.7534
	MDIVL	19.7054	19.6015	20.4821	19.5972	19.6279	<b>16.4441</b>	<b>18.3642</b>	19.0829	23.2351	<b>14.4779</b>
	KonIQ	0.4833	0.5283	0.5403	0.5741	0.4794	0.5114	<b>0.4323</b>	0.4359	<b>0.4300</b>	<b>0.3876</b>
	CLIVE	17.6477	19.4601	19.2687	19.6410	17.5379	19.2545	<b>17.3379</b>	17.6704	<b>17.4226</b>	<b>17.1298</b>
	CID2013	16.7826	22.5312	22.5098	20.2875	20.4756	20.9148	<b>15.5695</b>	<b>14.2467</b>	17.4576	<b>11.0931</b>
	SPAQ	20.1607	18.6684	20.2439	20.7546	<b>16.1107</b>	20.3706	17.4992	16.4737	<b>15.6132</b>	<b>14.5551</b>
	$AVG_D$	9.9112	11.4189	11.8547	11.7607	9.6716	10.2893	<b>9.2957</b>	<b>9.1318</b>	11.8545	<b>8.4532</b>
	$AVG_W$	7.7271	7.5637	8.0693	8.2122	<b>6.5589</b>	7.8258	6.8883	<b>6.6031</b>	6.8263	<b>5.9160</b>

is defined as the average of the two covariances. Therefore, the quality score is calculated as:

$$S(\vec{\mu}_d, \vec{\mu}_r, \sum_d, \sum_r) = \sqrt{\left( (\vec{\mu}_d - \vec{\mu}_r)^T \left( \frac{\sum_d + \sum_r}{2} \right)^{-1} (\vec{\mu}_d - \vec{\mu}_r) \right)}, \quad (8)$$

where  $\{\vec{\mu}_d, \sum_d\}$  and  $\{\vec{\mu}_r, \sum_r\}$  denote the mean vectors and covariance matrices of the estimated MVG models of the distorted image and the high-quality image set, respectively.

#### IV. EXPERIMENT

In this section, we first describe the basic experimental protocol to ensure a fair comparison with existing BIQA methods. We then evaluate the performance of the proposed model compared with classical and state-of-the-art models. Finally, we conduct extensive ablation studies to further evaluate the effectiveness of the custom modules and analyze the limitation of the proposed MDFS.

##### A. Experiment Protocol

1) *Training Dataset:* In this work, we collect a training dataset including 500 high-quality natural images with various

TABLE II: SROCC and KROCC comparisons of various IQA models under different distortion types on the TID2013 dataset.

Distortions	SROCC $\uparrow$									
	NIQE	QAC	PIQE	LPSI	ILNIQE	dipIQ	SNP-NIQE	NPQI	ContentSep	MDFS (Ours)
ANG	0.8187	0.7427	0.8555	0.7692	<b>0.8767</b>	<b>0.8653</b>	<b>0.8855</b>	0.6257	0.7997	0.8499
NCC	0.6701	0.7184	<b>0.7582</b>	0.4952	<b>0.8159</b>	<b>0.7687</b>	0.7323	0.2966	0.7341	0.7380
SCN	0.6659	0.1695	0.3354	<b>0.6967</b>	<b>0.9233</b>	0.5804	0.6507	0.0119	0.5806	<b>0.8150</b>
MN	<b>0.7464</b>	0.5927	0.5752	0.0468	0.5135	<b>0.7250</b>	<b>0.7383</b>	0.6624	0.6582	0.6490
HFN	0.8454	0.8628	<b>0.8923</b>	<b>0.9250</b>	0.8691	0.8642	0.8730	0.8214	0.8794	<b>0.8875</b>
IN	0.7437	<b>0.8003</b>	0.6901	0.4324	0.7556	<b>0.7878</b>	<b>0.8006</b>	0.5677	0.7138	0.7703
QN	0.8503	0.7089	0.7508	0.8536	<b>0.8714</b>	0.7991	<b>0.8573</b>	0.7732	0.7357	<b>0.8729</b>
GB	0.7969	0.8464	0.8280	0.8357	0.8145	<b>0.9046</b>	<b>0.8628</b>	0.7595	0.7852	<b>0.8614</b>
ID	0.5901	0.3381	<b>0.6442</b>	0.2487	<b>0.7494</b>	0.0690	0.6118	0.6403	0.5854	<b>0.8752</b>
JPEG	0.8427	0.8369	0.7929	<b>0.9122</b>	0.8343	<b>0.9115</b>	0.8775	0.8474	0.6507	<b>0.8952</b>
J2K	0.8890	0.7895	0.8536	<b>0.8983</b>	0.8583	<b>0.9194</b>	0.8813	0.8507	0.8242	<b>0.9326</b>
JTE	0.0727	0.0491	0.2287	0.0912	<b>0.3628</b>	<b>0.7085</b>	0.3214	0.0343	0.2019	<b>0.4233</b>
J2KTE	0.5250	0.4061	0.1129	<b>0.6106</b>	<b>0.6085</b>	0.3651	<b>0.6107</b>	0.0096	0.0962	0.5262
NEPN	<b>0.0687</b>	0.0479	0.0100	0.0522	<b>0.0809</b>	<b>0.3714</b>	0.0073	0.0621	0.0126	0.0165
LBD	0.1305	<b>0.2473</b>	0.1778	0.1374	0.1317	<b>0.2912</b>	0.0328	0.0901	<b>0.2882</b>	0.0889
MS	0.1627	<b>0.3060</b>	<b>0.2784</b>	<b>0.3406</b>	0.1843	0.0987	0.0649	0.0956	0.1191	0.1812
CC	0.0172	0.2067	0.0715	0.1994	0.0144	0.1369	0.4623	<b>0.2501</b>	<b>0.2840</b>	<b>0.2840</b>
CCS	0.2462	<b>0.3691</b>	0.2682	0.3017	0.1654	0.0700	0.1316	<b>0.3781</b>	0.1349	<b>0.5557</b>
MGN	0.6934	<b>0.7902</b>	0.7322	0.6960	0.6936	<b>0.7882</b>	0.7406	0.3958	<b>0.8026</b>	0.7590
CN	0.1914	0.1523	0.1475	0.0180	<b>0.3941</b>	<b>0.3909</b>	0.2242	0.1370	0.3119	<b>0.3181</b>
LCN	0.8025	0.6399	0.6369	0.2356	0.8287	<b>0.8513</b>	<b>0.8307</b>	0.3429	0.7700	<b>0.8458</b>
ICQ	0.7827	<b>0.8733</b>	<b>0.8119</b>	<b>0.8969</b>	0.7496	0.7562	0.7890	0.7556	0.0911	0.8043
CA	0.5620	0.6250	0.6756	<b>0.6953</b>	0.6793	<b>0.6998</b>	0.6339	0.5816	0.5206	<b>0.7177</b>
SSR	0.8340	0.7857	0.8229	<b>0.8580</b>	<b>0.8643</b>	0.7610	0.8284	0.8251	0.7245	<b>0.9196</b>

image sizes and content types from the DIV2K [44] dataset, which is non-overlapped with the testing dataset. However, we study the impact of different training datasets on our proposed method in the subsequent ablation experiment section.

2) *Testing Datasets*: To comprehensively evaluate the performance of the proposed model, various IQA datasets have been used to conduct extensive experiments, including CLIVE [45], CID 2013 [46], KonIQ [47], SPAQ [48], LIVE [49], CSIQ [50], TID2013 [51], KADID [52], MDLIVE [53], and MDIVL [54]. Specifically, the above datasets can be classified into two categories according to the degradation methods, where the first four and last six datasets belong to the realistic distortion-based datasets and synthetic distortion-based datasets, respectively.

3) *Evaluation Criteria*: Following [21], [55], the predicted objective scores are first mapped to the subjective scores by a nonlinear regression function as:

$$Q_i(x_i) = \gamma_1 \left( \frac{1}{2} - \frac{1}{1 + e^{\gamma_2(x_i - \gamma_3)}} \right) + \gamma_4 x_i + \gamma_5, \quad (9)$$

where  $x_i$  and  $Q_i$  are the predicted score and corresponding mapped score, respectively. The five parameters  $\gamma_1, \gamma_2, \gamma_3, \gamma_4$  and  $\gamma_5$  are determined by fitting the regression model. To evaluate the effectiveness of the proposed MDFS and comparison models, we employ four common evaluation criteria: Spearman rank order correlation coefficient (SROCC), Kendall rank order correlation coefficient (KROCC), Pearson linear correlation coefficient (PLCC), and Root-mean-square error (RMSE). It should be noted that a better IQA method should yield higher values for SROCC, KRCC, PLCC, and lower values for RMSE.

### B. Comparisons with State-of-the-Arts

In this subsection, we first compare the proposed MDFS with nine classical and state-of-the-art OU-BIQA methods, including NIQE [11], QAC [24], PIQE [12], LPSI [25], ILNIQE [21], dipIQ [35], SNP-NIQE [22], NPQI [23], and ContentSep [37]. Our focus is on demonstrating the advantages of our proposed MDFS across multiple datasets and diverse types of distortions. Detailed findings are presented below.

1) *Performance on Multiple IQA datasets*: Table I presents the performance comparison of various BIQA methods on ten datasets. Notably, the top three performers for each measurement criterion (*i.e.*, PLCC, SROCC, KROCC, and RMSE) are highlighted, with the first-ranked, second-ranked, and third-ranked IQA models emphasized in bold red, blue, and black, respectively. From the results, we can observe that the proposed MDFS significantly outperforms state-of-the-art OU-BIQA models on TID2013, KADID, MDLIVE, KonIQ, CID2013, and SPAQ datasets in terms of the PLCC, SROCC, KRCC, and RMSE. Additionally, on the LIVE, CSIQ, and CLIVE datasets, MDFS ranks top third in overall performance and is almost comparable with the state-of-the-art models.

To evaluate the overall performance across multiple datasets, two average measurements are applied as suggested in [56],

$$\bar{c} = \frac{\sum_{i=1}^N (c_i \cdot w_i)}{\sum_{i=1}^N w_i}, \quad (10)$$

where  $N$  denotes the number of datasets,  $c_i$  and  $w_i$  indicate the value of the measurement criteria and corresponding weight on the  $i$ -th dataset. Herein, we first set all the weights to 1 to obtain the results of the *Direct Average* ( $AVG_D$ ). Afterward,  $w_i$  is set to the number of the distorted images in the  $i$ -th dataset to obtain the results of the *Weighted Average* ( $AVG_W$ ). Due to the large difference in the range of quality

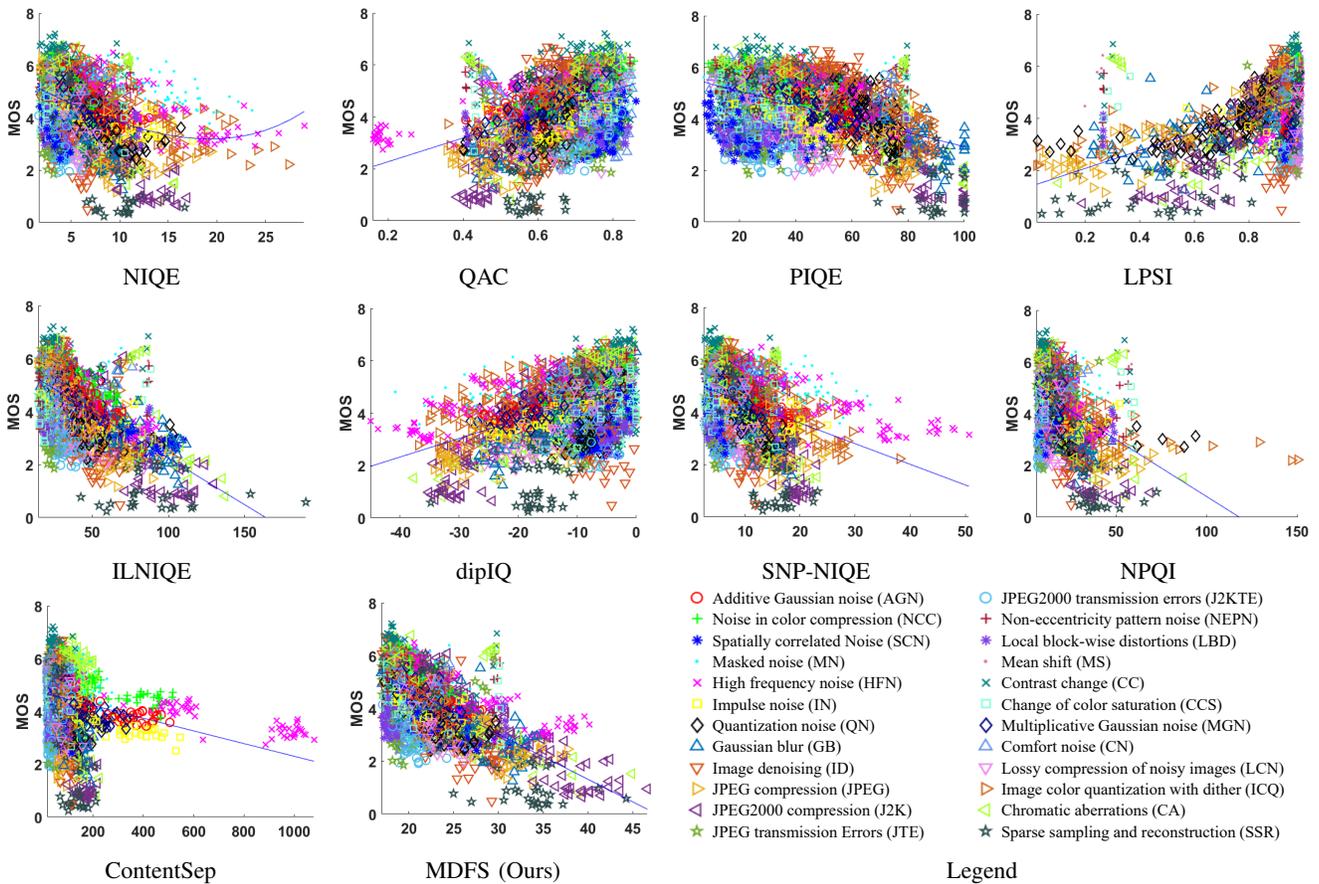


Fig. 2: Scatter plots of the mean opinion scores (MOS) versus the objective scores computed by the IQA models: (a) NIQE; (b) QAC; (c) PIQE; (d) LPSI; (e) ILNIQE; (f) dipIQ; (g) SNP-NIQE; (h) NPQI; (i) ContentSep, and (j) MDFS, respectively.

scores obtained by various IQA models, the average RMSE is not applicable for fairness. From Table I, the proposed MDFS is superior to all existing BIQA models in terms of both direct average and weighted average comparisons, which demonstrates the generality of the proposed model.

2) *Performance on Diverse Distortion Types:* In this section, we conduct a detailed and comprehensive experiment to evaluate the performance of the proposed MDFS model across various types of distortions using the TID2013 dataset. From Table II, one can observe that our model yields the most top-one and top-three performances compared with other IQA models in terms of SROCC. Specifically, in the comparisons in terms of SROCC, the proposed MDFS stands in the top three 14 times and ranks first 6 times, which is the same as the dipIQ model. The following are the ILNIQE, and LPSI, which ranked the top three 10 times and 9 times respectively, and ranked the top one 3 times and 4 times, respectively.

To visually compare the performance of BIQA models, Figure 2 provides scatter plots of subjective scores (*i.e.*, MOS) and objective scores, where all the objective scores are generated by the IQA models and further mapped using Equ. (9). The blue lines in each sub-figure represents the fitted line obtained from Equ. (9), indicating the “mean” value of the performance. For each distortion type, the predicted quality score (along the horizontal axis) is expected to be

close to the MOS value (along the vertical axis). Therefore, the closer the points representing a specific type of distortion are clustered around this line, the better the algorithm performs on that type of distortion. For example, for the JPEG2000 compression (J2K) distortion type, the predicted quality scores are closer to the corresponding fitted blue line in the scatter plot of the proposed MDFS, while predicted quality scores of dipIQ and LPSI are far away from the corresponding fitted blue line. This demonstrates that the quality scores predicted by the proposed MDFS are more consistent with the HVS regarding the J2K distortion, aligning with the results shown in Table II. Beyond the localized examination of scatter plots, an overarching perspective facilitates the observation of the interplay between various distortion types and the reference blue line. For instance, it is obvious that the proposed MDFS framework exhibits a closer fit to the blue line on various distortion types than other BIQA models. This coherent alignment substantiates the comprehensive superiority of the performance rendered by the proposed model, which is consistent with the results shown in Table I.

### C. Generalizability of MDFS

In this subsection, we evaluate the generalizability of the proposed MDFS model on four target-specific IQA datasets: underwater image (UWI), artificial intelligence generated con-

TABLE III: Performance comparisons of the original/retrained OU-BIQA models on four target-specific IQA datasets.

Datasets		SROCC $\uparrow$						KROCC $\uparrow$					
		NIQE		ILNIQE		MDFS (ours)		NIQE		ILNIQE		MDFS (ours)	
		Original	Retrained	Original	Retrained	Original	Retrained	Original	Retrained	Original	Retrained	Original	Retrained
UWI	SAUD	0.0616	0.1058	0.2818	0.2809	0.3353	<b>0.4431</b>	0.0417	0.0717	0.1920	0.1920	0.2285	<b>0.3098</b>
	UWIQA	0.4484	0.4803	0.4718	0.3440	0.3020	<b>0.6014</b>	0.3348	0.3576	0.3490	0.2473	0.2202	<b>0.4594</b>
AIGC	AGIQA	0.5338	-	0.5943	-	<b>0.6724</b>	-	0.3651	-	0.4194	-	<b>0.4816</b>	-
	AIGCIQA	0.5062	-	0.5692	-	<b>0.6992</b>	-	0.3422	-	0.3852	-	<b>0.4887</b>	-
Bird	MMQA-Birds	0.3098	0.0989	0.1226	0.2367	0.1883	<b>0.4974</b>	0.2119	0.0674	0.0842	0.1627	0.1273	<b>0.3473</b>
Face	GFIQA-20k	0.5011	0.6655	0.7142	0.7656	0.8331	<b>0.8359</b>	0.3495	0.4793	0.5184	0.5720	0.6410	<b>0.6479</b>
Datasets		PLCC $\uparrow$						RMSE $\downarrow$					
		NIQE		ILNIQE		MDFS (ours)		NIQE		ILNIQE		MDFS (ours)	
		Original	Retrained	Original	Retrained	Original	Retrained	Original	Retrained	Original	Retrained	Original	Retrained
UWI	SAUD	0.0876	0.1569	0.3141	0.3327	0.3663	<b>0.5334</b>	1.5500	1.5367	1.4772	1.4674	1.4479	<b>1.3162</b>
	UWIQA	0.4451	0.3373	0.4062	0.3368	0.4242	<b>0.5588</b>	0.1364	0.1434	0.1392	0.1434	0.1955	<b>0.1263</b>
AIGC	AGIQA	0.5391	-	0.6229	-	<b>0.6762</b>	-	0.8377	-	0.7805	-	<b>0.7350</b>	-
	AIGCIQA	0.5219	-	0.5641	-	<b>0.7047</b>	-	7.9454	-	7.6914	-	<b>6.6087</b>	-
Bird	MMQA-Birds	0.3125	0.1167	0.1492	0.2476	0.2260	<b>0.4531</b>	1.5988	1.6716	1.6642	1.6306	1.6395	<b>1.5004</b>
Face	GFIQA-20k	0.5006	0.6639	0.7022	0.6910	<b>0.8070</b>	0.7948	0.1580	0.1341	0.1299	0.1297	<b>0.1078</b>	0.1089

TABLE IV: Performance comparisons between our proposed OU-BIQA and different OA-BIQA methods on the public datasets.

	PaQ2PiQ	HyperIQA	MANIQA	VCRNet	MUSIQ	MDFS (Ours)	PaQ2PiQ	HyperIQA	MANIQA	VCRNet	MUSIQ	MDFS (Ours)	
SROCC	LIVE	0.4794	<b>0.7551</b>	<b>0.7793</b>	-	0.7335	<b>0.9361</b>	0.3557	<b>0.5515</b>	<b>0.5736</b>	-	0.5422	<b>0.7709</b>
	CSIQ	0.5643	0.5814	<b>0.6624</b>	<b>0.6806</b>	0.5878	<b>0.7774</b>	0.4000	0.4004	<b>0.4679</b>	<b>0.5031</b>	0.4129	<b>0.5823</b>
	TID2013	0.4011	0.3839	0.4510	<b>0.5116</b>	<b>0.4738</b>	<b>0.5363</b>	0.2838	0.2606	0.3175	<b>0.3667</b>	<b>0.3314</b>	<b>0.3824</b>
	KADID	0.3828	<b>0.4679</b>	0.4381	0.4443	<b>0.4640</b>	<b>0.5983</b>	0.2678	<b>0.3210</b>	0.3056	0.3099	<b>0.3230</b>	<b>0.4238</b>
	MDLIVE	<b>0.7737</b>	0.6594	0.5667	<b>0.8293</b>	<b>0.8647</b>	0.7579	<b>0.5726</b>	0.4738	0.4026	<b>0.6391</b>	<b>0.6734</b>	0.5623
	MDIVL	0.5361	<b>0.6174</b>	0.5110	0.4750	<b>0.5917</b>	<b>0.7890</b>	0.3777	<b>0.4378</b>	0.3594	0.3243	<b>0.4180</b>	<b>0.5911</b>
	KonIQ	<b>0.7213</b>	-	-	<b>0.6062</b>	-	<b>0.7333</b>	<b>0.5260</b>	-	-	<b>0.4268</b>	-	<b>0.5344</b>
	CLIVE	<b>0.7178</b>	<b>0.7612</b>	<b>0.8399</b>	0.5568	0.7216	0.4821	0.5293	<b>0.5612</b>	<b>0.6482</b>	0.3872	<b>0.5302</b>	0.3274
	CID2013	<b>0.8243</b>	0.7219	<b>0.8457</b>	0.5640	0.7685	<b>0.8571</b>	<b>0.6343</b>	0.5385	<b>0.6555</b>	0.3988	0.5784	<b>0.6706</b>
	SPAQ	0.6128	<b>0.8215</b>	0.0699	<b>0.7548</b>	<b>0.8323</b>	0.7408	0.4194	<b>0.6112</b>	0.0439	<b>0.5433</b>	<b>0.6227</b>	0.5347
$AV_{GD}$	0.6014	<b>0.6411</b>	0.5738	0.6025	<b>0.6709</b>	<b>0.7208</b>	0.4367	<b>0.4618</b>	0.4194	0.4332	<b>0.4925</b>	<b>0.5380</b>	
$AV_{GW}$	<b>0.5669</b>	0.4666	0.2516	<b>0.5868</b>	0.4766	<b>0.6854</b>	<b>0.4022</b>	0.3365	0.1779	<b>0.4171</b>	0.3474	<b>0.4966</b>	
PLCC	LIVE	0.4588	<b>0.7375</b>	<b>0.7623</b>	-	0.6722	<b>0.8558</b>	24.2770	<b>18.4522</b>	<b>17.6835</b>	-	20.2281	<b>14.1344</b>
	CSIQ	0.6360	0.5609	<b>0.6549</b>	<b>0.7514</b>	0.6276	<b>0.7907</b>	0.2026	0.2174	<b>0.1999</b>	<b>0.1732</b>	0.2044	<b>0.1607</b>
	TID2013	<b>0.5776</b>	0.4427	0.4873	<b>0.6215</b>	0.5771	<b>0.6242</b>	<b>1.0120</b>	1.1116	1.0830	<b>0.9711</b>	1.0124	<b>0.9685</b>
	KADID	0.4356	<b>0.4919</b>	0.4788	0.4819	<b>0.5035</b>	<b>0.5939</b>	0.9745	<b>0.9426</b>	0.9513	0.9486	<b>0.9354</b>	<b>0.8710</b>
	MDLIVE	0.8160	0.7554	0.6609	<b>0.8582</b>	<b>0.8824</b>	<b>0.8226</b>	10.9323	12.3933	14.1928	<b>9.7094</b>	<b>8.8976</b>	<b>10.7534</b>
	MDIVL	0.5223	<b>0.6246</b>	0.5342	0.4881	<b>0.5798</b>	<b>0.7953</b>	20.3651	<b>18.6496</b>	20.1880	20.8439	<b>19.4578</b>	<b>14.4779</b>
	KonIQ	<b>0.7259</b>	-	-	<b>0.6231</b>	-	<b>0.7123</b>	<b>0.3798</b>	-	-	<b>0.4319</b>	-	<b>0.3876</b>
	CLIVE	<b>0.7706</b>	<b>0.7739</b>	<b>0.8481</b>	0.5655	0.7498	0.5364	<b>12.9364</b>	<b>12.8538</b>	<b>10.7528</b>	16.7398	13.4291	17.1298
	CID2013	<b>0.8261</b>	0.7856	<b>0.8326</b>	0.6077	0.8224	<b>0.8717</b>	<b>12.7582</b>	14.0090	<b>12.5406</b>	17.9808	12.8806	<b>11.0931</b>
	SPAQ	0.5663	<b>0.8301</b>	0.0666	<b>0.7656</b>	<b>0.8272</b>	0.7177	17.2262	<b>11.6554</b>	20.8550	<b>13.4464</b>	<b>11.7459</b>	14.5551
$AV_{GD}$	0.6335	<b>0.6670</b>	0.5917	0.6403	<b>0.6936</b>	<b>0.7321</b>	10.1064	10.0317	10.9385	<b>9.0272</b>	<b>9.8657</b>	<b>8.4532</b>	
$AV_{GW}$	<b>0.5852</b>	0.4815	0.2652	<b>0.6155</b>	0.4946	<b>0.6803</b>	6.9253	<b>5.1090</b>	7.7004	<b>5.5300</b>	<b>5.1392</b>	5.9160	

tent (AIGC), birds images (Bird), and human face images (Face). We compare MDFS with two well-established models, NIQE and ILNIQE, since they have relatively good performance and can be easily retrained with new high-quality images. The experiment includes two UWI datasets: SAUD [57] and UWIQA [58], two AIGC datasets: AGIQA [59] and AIGCIQA [60], one Bird dataset: MMQA-Birds [61], and one Face dataset: GFIQA-20k [62] for testing. Specifically, for high-quality image dataset for retraining, we select the first 200 high-quality enhanced images from UID2021 [63] as the high-quality UWI images, 200 high-quality bird images

from the Caltech-UCSD Birds-200-2011 dataset [64], and 200 high-quality face images from the CelebA-HQ [65] dataset. Notably, the high-quality image dataset for retraining has a distribution similar to the corresponding testing dataset, but there is no overlap between the images used for training and testing. This ensures that the performance of the model is evaluated on unseen data, testing its ability to generalize to new and unobserved images.

From the results presented in Table III, one can observe that the retrained model generally outperforms the original one in most cases. This improvement can be attributed to the fact that

KADID	NIQE	QAC	PIQE	LPSI	ILNIQE	dipIQ	SNP-NIQE	NPQI	ContentSep	MDFS (Ours)
NIQE	0	0	1	0	1	1	1	1	1	0
QAC	1	0	1	0	1	1	1	1	1	0
PIQE	0	0	0	0	1	1	0	0	1	0
LPSI	1	0	1	0	1	1	1	1	1	0
ILNIQE	0	0	0	0	0	1	0	0	1	0
dipIQ	0	0	0	0	0	0	0	0	1	0
SNP-NIQE	0	0	1	0	1	1	0	1	1	0
NPQI	0	0	1	0	1	1	0	0	1	0
ContentSep	0	0	0	0	0	0	0	0	0	0
MDFS (Ours)	1	1	1	1	1	1	1	1	1	0

(a)

TID2013	NIQE	QAC	PIQE	LPSI	ILNIQE	dipIQ	SNP-NIQE	NPQI	ContentSep	MDFS (Ours)
NIQE	0	0	1	0	1	1	1	1	1	0
QAC	1	0	1	0	1	1	1	1	1	0
PIQE	0	0	0	0	1	1	0	0	1	0
LPSI	1	0	1	0	1	1	1	1	1	0
ILNIQE	0	0	0	0	0	1	0	0	1	0
dipIQ	0	0	0	0	0	0	0	0	1	0
SNP-NIQE	0	0	1	0	1	1	0	1	1	0
NPQI	0	0	1	0	1	1	0	0	1	0
ContentSep	0	0	0	0	0	0	0	0	0	0
MDFS (Ours)	1	1	1	1	1	1	1	1	1	0

(b)

CSIQ	NIQE	QAC	PIQE	LPSI	ILNIQE	dipIQ	SNP-NIQE	NPQI	ContentSep	MDFS (Ours)
NIQE	0	0	1	0	1	1	1	1	1	1
QAC	1	0	1	0	1	1	1	1	1	1
PIQE	0	0	0	0	1	1	0	0	1	0
LPSI	1	1	1	0	1	1	1	1	1	1
ILNIQE	0	0	0	0	0	1	0	0	1	0
dipIQ	0	0	0	0	0	0	0	0	1	0
SNP-NIQE	0	0	1	0	1	1	0	1	1	1
NPQI	0	0	1	0	1	1	0	0	1	0
ContentSep	0	0	0	0	0	0	0	0	0	0
MDFS (Ours)	0	0	1	0	1	1	0	1	1	0

(c)

Fig. 3: The statistically significant test results of various OU-BIQA methods on the (a) KADID, (b) TID2013, and (c) CSIQ datasets. A value of “1” indicates that the model in the row is significantly better than the model in the column.

the retrained model is learned from target-specific images that are close to the testing images, whereas the original model is trained on general high-quality images. In comparison to the noticeable improvement in performance after retraining on two tasks (i.e., bird and underwater datasets), the results of the model retrained on the face dataset remain remarkably close to the original results. We speculate that this close resemblance may be attributed to the saturation of model performance on this task. Additionally, the retrained MDFS model demonstrates superior performance compared to other methods across all three image categories, showcasing the effectiveness of the proposed MDFS framework in generalizing to diverse target-specific IQA tasks.

#### D. Significant Test

In this subsection, the F-test is adopted as statistical analysis to illustrate the superiority of the proposed MDFS model compared with other OU-BIQA models. The results of the F-test between each pair of BIQA models on KADID, TID2013, and CSIQ datasets are presented in Figure 3. A value of “1” (highlighted in yellow) indicates that the model in the row significantly outperforms the model in the column, while a value of “0” (highlighted in green) indicates that both algorithms in the column and row are statistically equivalent. The results depicted in Figure 3 reveal that the proposed MDFS outperforms all the classical and state-of-the-art OU-BIQA models on both the KADID and TID2013 datasets. On the CSIQ dataset, MDFS demonstrates superiority over most of the comparison OU-IQA methods (i.e., PIQE, ILNIQE, dipIQ, NPQI, and ContentSep). These findings highlight the statistical significance of the performance advantage of MDFS on different datasets.

#### E. Cross-dataset Comparison with OA-BIQA Methods

In this subsection, we conduct a comprehensive comparison of the proposed MDFS with five OA-BIQA methods, including PaQ2PiQ [66], HyperIQA [31], MANIQA [34], VCRNet [67], and MUSIQ [68]. From the results in Table IV, one can observe that MDFS outperforms the compared OA-BIQA methods on ten public datasets in most cases. Furthermore,

TABLE V: Computational efficiency compared with state-of-the-art methods on the CID2013 and TID2013 datasets.

Methods	Programming Language	CID2013	TID2013
NIQE	MATLAB	0.2216	<b>0.0263</b>
QAC	MATLAB	0.4259	0.0618
PIQE	MATLAB	0.3502	0.0413
LPSI	MATLAB	<b>0.1304</b>	<b>0.0194</b>
ILNIQE	MATLAB	1.8163	1.5908
dipIQ	MATLAB	1.8800	0.9905
SNP-NIQE	MATLAB	11.7991	0.8245
NPQI	MATLAB	7.8984	0.8534
ContentSep	Python	<b>0.0470</b>	0.1518
MDFS (Ours)	Python	<b>0.0558</b>	<b>0.0289</b>

both the direct average score and the weighted average score of the proposed MDFS demonstrate superior performance compared to all OA-BIQA methods in terms of SROCC, KROCC, and PLCC. This indicates that MDFS exhibits remarkable performance compared to algorithms that require training based on subjective scores. These findings underscore the excellence of the proposed MDFS approach and its suitability for a wide range of IQA applications. We believe this is mainly due to the fact that our algorithm takes advantage of the visual model pre-trained on a large data set, which enables it to extract more general features and conduct statistical analysis.

#### F. Computational Efficiency

In this subsection, we conduct the computational efficiency comparison experiment, which is an important factor in evaluating the performance of IQA models in real-world applications. The computational efficiency of all BIQA models is assessed using the CID2013 (474 distorted images with a resolution of  $1600 \times 1200$ ) and TID2013 (3000 distorted images with a resolution of  $512 \times 384$ ) datasets, which have varying image resolutions, to compute the average running time per image. All the experiments run on the computer with an Intel i7-9700K CPU @ 3.60GHz and an NVIDIA GeForce RTX 2080 Ti GPU. Furthermore, all the codes of IQA models are performed under the same suggestion of the corresponding authors. The results in Table V indicate that the proposed MDFS model demonstrates similar computational speed to

TABLE VI: Ablation studies on various training datasets. The top three are marked in bold red, blue, and black, respectively.

Training dataset	Criteria	LIVE	CSIQ	TID2013	KADID	MDLIVE	MDIVL	KonIQ	CLIVE	CID2013	SPAQ
Waterloo	SROCC	<b>0.9072</b>	<b>0.7964</b>	<b>0.5752</b>	<b>0.6214</b>	<b>0.7763</b>	<b>0.8331</b>	0.4542	0.3998	<b>0.6249</b>	<b>0.7191</b>
	KROCC	<b>0.7296</b>	<b>0.5998</b>	<b>0.4130</b>	<b>0.4431</b>	<b>0.5644</b>	<b>0.6332</b>	0.3114	0.2705	<b>0.4557</b>	<b>0.5195</b>
	PLCC	0.7925	<b>0.7980</b>	<b>0.6380</b>	<b>0.6434</b>	<b>0.7212</b>	<b>0.7392</b>	0.4687	0.4659	<b>0.5631</b>	<b>0.6605</b>
	RMSE	16.6633	<b>0.1582</b>	<b>0.9546</b>	<b>0.8288</b>	<b>13.1012</b>	<b>16.0835</b>	0.4878	17.9591	<b>18.7101</b>	<b>15.6932</b>
D-NIQE	SROCC	0.8900	<b>0.7792</b>	<b>0.5524</b>	<b>0.6299</b>	0.7245	0.6629	<b>0.4721</b>	<b>0.4404</b>	0.4894	0.6897
	KROCC	0.7052	<b>0.5856</b>	<b>0.3965</b>	<b>0.4509</b>	0.5183	0.4705	<b>0.3248</b>	<b>0.3012</b>	0.3504	0.4926
	PLCC	<b>0.7998</b>	0.7859	<b>0.6266</b>	<b>0.6541</b>	0.6189	0.5388	<b>0.4709</b>	<b>0.4711</b>	0.4366	0.6418
	RMSE	<b>16.3997</b>	0.1623	<b>0.9662</b>	<b>0.8189</b>	14.8549	20.1192	<b>0.4871</b>	<b>17.9028</b>	20.3681	16.0279
KADID	SROCC	<b>0.9093</b>	<b>0.8101</b>	<b>0.5507</b>	-	<b>0.7597</b>	<b>0.7352</b>	<b>0.5535</b>	<b>0.5313</b>	<b>0.5996</b>	<b>0.7387</b>
	KROCC	<b>0.7331</b>	<b>0.6165</b>	<b>0.3937</b>	-	<b>0.5550</b>	<b>0.5351</b>	<b>0.3893</b>	<b>0.3666</b>	<b>0.4428</b>	<b>0.5374</b>
	PLCC	<b>0.9127</b>	<b>0.8161</b>	<b>0.6327</b>	-	<b>0.7126</b>	<b>0.7264</b>	<b>0.4949</b>	<b>0.5456</b>	<b>0.6432</b>	<b>0.6821</b>
	RMSE	<b>11.1664</b>	<b>0.1519</b>	<b>0.9600</b>	-	<b>13.2674</b>	<b>16.4126</b>	<b>0.4798</b>	<b>17.0093</b>	<b>17.3358</b>	<b>15.2846</b>
DIV2K (Ours)	SROCC	<b>0.9361</b>	0.7774	0.5363	<b>0.5983</b>	<b>0.7579</b>	<b>0.7890</b>	<b>0.7333</b>	<b>0.4821</b>	<b>0.8571</b>	<b>0.7408</b>
	KROCC	<b>0.7709</b>	0.5823	0.3824	<b>0.4238</b>	<b>0.5623</b>	<b>0.5911</b>	<b>0.5344</b>	<b>0.3274</b>	<b>0.6706</b>	<b>0.5347</b>
	PLCC	<b>0.8558</b>	<b>0.7907</b>	0.6242	<b>0.5939</b>	<b>0.8226</b>	<b>0.7953</b>	<b>0.7123</b>	<b>0.5364</b>	<b>0.8717</b>	<b>0.7177</b>
	RMSE	<b>14.1344</b>	<b>0.1607</b>	0.9685	<b>0.8710</b>	<b>10.7534</b>	<b>14.4779</b>	<b>0.3876</b>	<b>17.1298</b>	<b>11.0931</b>	<b>14.5551</b>

TABLE VII: Ablation studies on various network backbones.

Datasets	Criteria	ConvNet	Inception	PNAS	ResNet	VGG	ViT	EN
LIVE	SROCC	0.8138	0.4559	0.8871	<b>0.8872</b>	<b>0.8896</b>	0.4675	<b>0.9361</b>
	KROCC	0.5991	0.3206	<b>0.6995</b>	0.6887	<b>0.6943</b>	0.3193	<b>0.7709</b>
	RMSE	15.7004	23.0483	16.3517	<b>12.3457</b>	<b>12.1704</b>	24.2031	<b>14.1344</b>
	PLCC	0.8184	0.5370	0.8011	<b>0.8921</b>	<b>0.8953</b>	0.464	<b>0.8558</b>
TID2013	SROCC	0.4106	0.2559	<b>0.5010</b>	<b>0.4495</b>	0.4160	0.3403	<b>0.5363</b>
	KROCC	0.2756	0.1738	<b>0.3491</b>	<b>0.3085</b>	0.2857	0.2318	<b>0.3824</b>
	RMSE	1.1332	1.1526	<b>1.0550</b>	1.0866	<b>1.0418</b>	1.1345	<b>0.9685</b>
	PLCC	0.4054	0.3682	<b>0.5251</b>	0.4813	<b>0.5420</b>	0.4031	<b>0.6242</b>

TABLE VIII: Ablation study on various window sizes.

Datasets	Criteria	$s_w=3$	$s_w=5$	$s_w=7$	Ours
CID2013	SROCC	0.8410	<b>0.8489</b>	<b>0.8495</b>	<b>0.8571</b>
	KROCC	0.6517	<b>0.6606</b>	<b>0.6617</b>	<b>0.6706</b>
	RMSE	<b>11.6844</b>	<b>13.7513</b>	13.7975	<b>11.0931</b>
	PLCC	<b>0.8565</b>	<b>0.7944</b>	0.7928	<b>0.8717</b>
LIVE	SROCC	0.9337	<b>0.9352</b>	<b>0.9346</b>	<b>0.9361</b>
	KROCC	0.7665	<b>0.7685</b>	<b>0.7669</b>	<b>0.7709</b>
	RMSE	<b>9.9014</b>	<b>12.8541</b>	<b>9.8584</b>	14.1344
	PLCC	<b>0.9320</b>	<b>0.8824</b>	<b>0.9326</b>	0.8558

the fastest algorithms across the two datasets. Moreover, the proposed model manifests pronounced advantages specifically in the context of high-resolution imagery.

### G. Ablation Studies

This subsection conducts extensive ablation studies to evaluate the impact of each component in the proposed MDIFS model, including training dataset, network backbones, window size, quality calculation, and contrast feature.

1) *Training Dataset*: In our study, we utilize the DIV2K dataset for training purposes due to its excellent image quality and diverse content. However, we also explored the performance of using other high-quality image datasets, specifically the Waterloo dataset, the reference images of KADID, and the training dataset of NIQE (referred to as D-NIQE). The evaluation results are summarized in Table VI. It is evident that training on the DIV2K dataset led to relatively better performance compared to training on the other datasets. Furthermore, models trained on those alternative datasets still delivered commendable results when compared to existing algorithms, as demonstrated in Table I. This highlights the robustness and versatility of the proposed MDIFS model.

2) *Network Backbone*: In our study, we explored the use of various existing networks as the backbone for the MDIFS to provide feature maps, including VGG [39], ResNet [41], PNAS [69], ConvNet [40], ViT [38], and EfficientNet [42]. Table VII presents the SROCC results on the LIVE and

TID2013 datasets with various datasets. These results indicate that utilizing alternative backbone networks, such as ResNet, VGG, PNAS, and ViT did not yield results as promising as when employing EfficientNet. Therefore, we use EfficientNet as the preferred backbone for the feature extraction module as it consistently shows the most promising results.

3) *Window Size*: In our proposed MDIFS model, we introduce a novel window size calculation method, as described in Equ. (4), which dynamically adjusts the window size based on the dimensions of the input image. To evaluate its effectiveness, we compare the proposed dynamic window size with three fixed window sizes: 3, 5, and 7. The results in Table VIII demonstrate that the proposed dynamic window algorithm not only enhances the accuracy of the proposed MDIFS model but also exhibits robustness across various input image sizes. This indicates that our approach adapts well to images of different dimensions, making it a versatile and effective solution.

4) *Contrast Feature*: In our proposed MDIFS model, we use the standard deviation as a measure of contrast to extract HVS-sensitive information. Herein, we investigate the performance of various contrast information extraction methods. To be specific, *w.o. std* represents the MDIFS model without weighting map (i.e., standard deviation); *w. var* refers to the MDIFS with variance as weighting map; *w. entropy* indicates the MDIFS with entropy feature as weighting map; and *w. w'* denotes

TABLE IX: Ablation study on different contrast features.

Datasets		w/o. std	w. var	w. entropy	w. $w'$	w. std (ours)
SROCC	LIVE	0.9121	<b>0.9136</b>	0.9124	<b>0.9136</b>	<b>0.9361</b>
	CSIQ	0.7321	<b>0.7333</b>	0.7308	<b>0.7333</b>	<b>0.7774</b>
	TID2013	0.5181	<b>0.5226</b>	0.5213	<b>0.5226</b>	<b>0.5363</b>
	KADID	0.5716	<b>0.5828</b>	0.5713	<b>0.5820</b>	<b>0.5983</b>
	MDLIVE	0.6702	<b>0.6902</b>	0.6761	<b>0.6871</b>	<b>0.7579</b>
	MDIVL	<b>0.8008</b>	0.7971	<b>0.7993</b>	<b>0.7976</b>	0.7890
	KonIQ	0.5778	<b>0.5842</b>	0.5796	<b>0.5838</b>	<b>0.7333</b>
	CLIVE	0.3596	<b>0.3655</b>	<b>0.3682</b>	0.3648	<b>0.4821</b>
	CID2013	<b>0.5405</b>	0.5370	0.5288	<b>0.5395</b>	<b>0.8571</b>
	SPAQ	0.5626	<b>0.6007</b>	0.5765	<b>0.5974</b>	<b>0.7408</b>

TABLE X: Ablation study on different distance calculation methods.

Datasets		MMD	EMD	SWD	KL	MDFS (Ours)
SROCC	LIVE	0.4710	<b>0.6531</b>	0.3808	<b>0.6497</b>	<b>0.9361</b>
	CSIQ	0.2427	<b>0.4740</b>	0.1930	<b>0.4662</b>	<b>0.7774</b>
	TID2013	<b>0.2583</b>	0.2379	0.0693	<b>0.2340</b>	<b>0.5363</b>
	KADID	0.2202	<b>0.2602</b>	0.1368	<b>0.2520</b>	<b>0.5983</b>
	MDLIVE	<b>0.3429</b>	0.3301	0.2052	<b>0.3343</b>	<b>0.7579</b>
	MDIVL	0.3455	<b>0.4510</b>	0.0007	<b>0.4516</b>	<b>0.7890</b>
	KonIQ	0.2458	<b>0.3239</b>	0.0963	<b>0.3182</b>	<b>0.7333</b>
	CLIVE	0.0501	<b>0.1873</b>	0.0205	<b>0.1791</b>	<b>0.4821</b>
	CID2013	0.0187	<b>0.1980</b>	<b>0.3871</b>	0.1952	<b>0.8571</b>
	SPAQ	0.3828	<b>0.5144</b>	0.3739	<b>0.5105</b>	<b>0.7408</b>

using a new weighting map  $w'$  defined as follows:

$$w' = 1/(1 + e^{-F_\sigma^{2\sigma}/(F_\sigma^\mu + \delta)}), \quad (11)$$

where  $F_\sigma^\mu$  and  $F_\sigma^\sigma$  represent the mean and standard deviation of the  $F_\sigma$ , respectively.  $\delta$  is a small positive number (e.g.,  $\delta = 1 \times e^{-12}$ ) to prevent the denominator from being zero. The results in Table IX demonstrate that using standard deviation in the MDFS yields superior performance compared to other variants. The potential reason may be that the standard deviation can better capture HVS-sensitive features.

5) *Quality Calculation*: We explored the performance of various distance algorithms to calculate the final quality score, including Maximum Mean Discrepancy (MMD) [70], Earth Mover’s distance (EMD) [71], Sliced Wasserstein Distance (SWD) [72], and Kullback-Leibler (KL) Divergence [73]. From the results in Table X, it is evident that employing the MVG distance yields superior and robust results compared with other distance algorithms.

#### H. Limitation

The proposed MDFS model undoubtedly presents significant advancements in the OU-BIQA task. However, it also comes with certain limitations. Firstly, the proposed method relies on the high-quality image dataset for the learning of reference features. To efficiently identify the difference between high-quality images and distorted images, both the high visual quality and content diversity are required for the images used for learning. Secondly, the proposed method calculates the quality score based on the statistical feature analysis of the overall image, therefore ignoring the rationality of the local patches. For instance, NEPN and LBD distortions

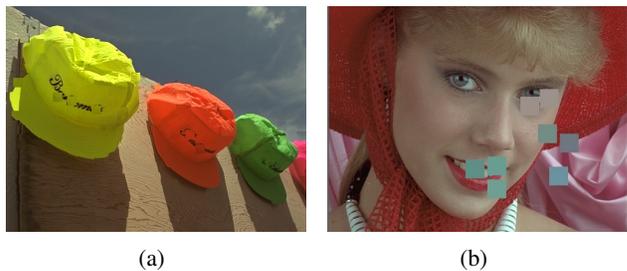


Fig. 4: Failure cases. The distorted images are generated by the (a) NEPN and (b) LBD distortions in TID2013.

alter local patch content but have minimal impact on overall statistical features, as illustrated in Figure 4. Consequently, MDFS demonstrates poor performance on images with these distortion types, as evidenced in Table II. In future work, it may be worthwhile to enhance the accuracy of the BIQA model on these distortion types by incorporating the local positional features of the images.

#### V. CONCLUSION

In this paper, we introduce a novel opinion-unaware blind image quality assessment (OU-BIQA) model called the Multi-scale Deep Feature Statistic (MDFS) model, which eliminates the need for human-rated data during training. The core idea of our approach involves integrating multi-scale deep features with a traditional statistical analysis model for OU-BIQA. On one hand, deep features provide richer and more expressive representations compared to conventional features. On the other hand, the statistical analysis model is highly efficient and stable, making the training process more cost-effective. Experimental results across various datasets demonstrate that our model achieves superior consistency with human visual perception compared to existing BIQA methods, while also exhibiting improved generalizability across diverse target-specific BIQA tasks.

While our research has made significant strides, it is imperative to acknowledge the limitation of some distortion types that do not significantly impact the global statistical data of the image. This results in a subpar performance of most IQA metrics on these distortion types. Future research endeavors could explore novel methodologies to address this limitation, potentially involving finer-grained distortion analysis based on local information or tailored processing strategies for different distortion types.

#### REFERENCES

- [1] Y. Li, S. Wang, X. Zhang, S. Wang, S. Ma, and Y. Wang, “Quality assessment of end-to-end learned image compression: The benchmark and objective measure,” in *Proceedings of the 29th ACM International Conference on Multimedia*, 2021, pp. 4297–4305.
- [2] C. Saharia, J. Ho, W. Chan, T. Salimans, D. J. Fleet, and M. Norouzi, “Image super-resolution via iterative refinement,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [3] Z. Ni, W. Yang, S. Wang, L. Ma, and S. Kwong, “Towards unsupervised deep image enhancement with generative adversarial network,” *IEEE Transactions on Image Processing*, vol. 29, pp. 9140–9151, 2020.

- [4] X. Wang, Q. Jiang, F. Shao, K. Gu, G. Zhai, and X. Yang, "Exploiting local degradation characteristics and global statistical properties for blind quality assessment of tone-mapped hdr images," *IEEE Transactions on Multimedia*, vol. 23, pp. 692–705, 2020.
- [5] Q. Li, W. Lin, J. Xu, and Y. Fang, "Blind image quality assessment using statistical structural and luminance features," *IEEE Transactions on Multimedia*, vol. 18, no. 12, pp. 2457–2469, 2016.
- [6] Y. Fang, J. Yan, R. Du, Y. Zuo, W. Wen, Y. Zeng, and L. Li, "Blind quality assessment for tone-mapped images by analysis of gradient and chromatic statistics," *IEEE Transactions on Multimedia*, vol. 23, pp. 955–966, 2020.
- [7] K. Ma, W. Liu, K. Zhang, Z. Duanmu, Z. Wang, and W. 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.
- [8] J. Guan, S. Yi, X. Zeng, W.-K. Cham, and X. Wang, "Visual importance and distortion guided deep image quality assessment framework," *IEEE Transactions on Multimedia*, vol. 19, no. 11, pp. 2505–2520, 2017.
- [9] X. Wang, J. Xiong, and W. Lin, "Visual interaction perceptual network for blind image quality assessment," *IEEE Transactions on Multimedia*, 2023.
- [10] C. Yang, X. Zhang, P. An, L. Shen, and C.-C. J. Kuo, "Blind image quality assessment based on multi-scale klt," *IEEE Transactions on Multimedia*, vol. 23, pp. 1557–1566, 2020.
- [11] A. Mittal, R. Soundararajan, and A. C. Bovik, "Making a "completely blind" image quality analyzer," *IEEE Signal Processing Letters*, vol. 20, no. 3, pp. 209–212, 2012.
- [12] N. Venkatanath, D. Praneeth, M. C. Bh, S. S. Channappayya, and S. S. Medasani, "Blind image quality evaluation using perception based features," in *Proceedings of the 21th IEEE National Conference on Communications*, 2015, pp. 1–6.
- [13] K. Ma, X. Liu, Y. Fang, and E. P. Simoncelli, "Blind image quality assessment by learning from multiple annotators," in *Proceedings of the IEEE International Conference on Image Processing*, 2019, pp. 2344–2348.
- [14] Z. Wang, Q. Jiang, S. Zhao, W. Feng, and W. Lin, "Deep blind image quality assessment powered by online hard example mining," *IEEE Transactions on Multimedia*, 2023.
- [15] K. Sim, J. Yang, W. Lu, and X. Gao, "Blind stereoscopic image quality evaluator based on binocular semantic and quality channels," *IEEE Transactions on Multimedia*, vol. 24, pp. 1389–1398, 2021.
- [16] Q. Jiang, W. Zhou, X. Chai, G. Yue, F. Shao, and Z. Chen, "A full-reference stereoscopic image quality measurement via hierarchical deep feature degradation fusion," *IEEE Transactions on Instrumentation and Measurement*, vol. 69, no. 12, pp. 9784–9796, 2020.
- [17] Q. Jiang, Z. Peng, F. Shao, K. Gu, Y. Zhang, W. Zhang, and W. Lin, "Stereoars: Quality evaluation for stereoscopic image retargeting with binocular inconsistency detection," *IEEE Transactions on Broadcasting*, vol. 68, no. 1, pp. 43–57, 2021.
- [18] Z. Peng, Q. Jiang, F. Shao, W. Gao, and W. Lin, "Lggd+: Image retargeting quality assessment by measuring local and global geometric distortions," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 6, pp. 3422–3437, 2021.
- [19] J. Xu, W. Zhou, H. Li, F. Li, and Q. Jiang, "Quality assessment of multi-exposure image fusion by synthesizing local and global intermediate references," *Displays*, vol. 74, p. 102188, 2022.
- [20] E. P. Simoncelli and B. A. Olshausen, "Natural image statistics and neural representation," *Annual Review of Neuroscience*, vol. 24, no. 1, pp. 1193–1216, 2001.
- [21] L. Zhang, L. Zhang, and A. C. Bovik, "A feature-enriched completely blind image quality evaluator," *IEEE Transactions on Image Processing*, vol. 24, no. 8, pp. 2579–2591, 2015.
- [22] Y. Liu, K. Gu, Y. Zhang, X. Li, G. Zhai, D. Zhao, and W. Gao, "Unsupervised blind image quality evaluation via statistical measurements of structure, naturalness, and perception," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 4, pp. 929–943, 2019.
- [23] Y. Liu, K. Gu, X. Li, and Y. Zhang, "Blind image quality assessment by natural scene statistics and perceptual characteristics," *ACM Transactions on Multimedia Computing, Communications, and Applications*, vol. 16, no. 3, pp. 1–91, 2020.
- [24] W. Xue, L. Zhang, and X. Mou, "Learning without human scores for blind image quality assessment," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 995–1002.
- [25] Q. Wu, Z. Wang, and H. Li, "A highly efficient method for blind image quality assessment," in *Proceedings of the IEEE International Conference on Image Processing*, 2015, pp. 339–343.
- [26] A. K. Moorthy and A. C. Bovik, "Blind image quality assessment: From natural scene statistics to perceptual quality," *IEEE Transactions on Image Processing*, vol. 20, no. 12, pp. 3350–3364, 2011.
- [27] M. A. Saad, A. C. Bovik, and C. Charrier, "Blind image quality assessment: A natural scene statistics approach in the dct domain," *IEEE Transactions on Image Processing*, vol. 21, no. 8, pp. 3339–3352, 2012.
- [28] Q. Wu, H. Li, F. Meng, K. N. Ngan, B. Luo, C. Huang, and B. Zeng, "Blind image quality assessment based on multichannel feature fusion and label transfer," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 26, no. 3, pp. 425–440, 2015.
- [29] S. Sun, T. Yu, J. Xu, W. Zhou, and Z. Chen, "Graphiqa: Learning distortion graph representations for blind image quality assessment," *IEEE Transactions on Multimedia*, 2022.
- [30] W. Zhang, D. Li, C. Ma, G. Zhai, X. Yang, and K. Ma, "Continual learning for blind image quality assessment," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [31] S. Su, Q. Yan, Y. Zhu, C. Zhang, X. Ge, J. Sun, and Y. Zhang, "Blindly assess image quality in the wild guided by a self-adaptive hyper network," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3667–3676.
- [32] Y. Zhu, Y. Li, W. Sun, X. Min, G. Zhai, and X. Yang, "Blind image quality assessment via cross-view consistency," *IEEE Transactions on Multimedia*, 2022.
- [33] J. Kim and S. Lee, "Fully deep blind image quality predictor," *IEEE Journal of Selected Topics in Signal Processing*, vol. 11, no. 1, pp. 206–220, 2016.
- [34] S. Yang, T. Wu, S. Shi, S. Lao, Y. Gong, M. Cao, J. Wang, and Y. Yang, "Maniqa: Multi-dimension attention network for no-reference image quality assessment," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 1191–1200.
- [35] K. Ma, W. Liu, T. Liu, Z. Wang, and D. Tao, "Dipiqa: Blind image quality assessment by learning-to-rank discriminative image pairs," *IEEE Transactions on Image Processing*, vol. 26, no. 8, pp. 3951–3964, 2017.
- [36] P. Chen, L. Li, Q. Wu, and J. Wu, "Spiqa: A self-supervised pre-trained model for image quality assessment," *IEEE Signal Processing Letters*, vol. 29, pp. 513–517, 2022.
- [37] N. C. Babu, V. Kannan, and R. Soundararajan, "No reference opinion unaware quality assessment of authentically distorted images," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 2459–2468.
- [38] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *Proceedings of the International Conference on Learning Representations*, 2021.
- [39] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proceedings of the International Conference on Learning Representations*, 2015.
- [40] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, "A convnet for the 2020s," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 11 976–11 986.
- [41] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [42] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *Proceedings of the PMLR International Conference on Machine Learning*, 2019, pp. 6105–6114.
- [43] Y. Fu, H. Zeng, L. Ma, Z. Ni, J. Zhu, and K.-K. Ma, "Screen content image quality assessment using multi-scale difference of gaussian," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 9, pp. 2428–2432, 2018.
- [44] E. Agustsson and R. Timofte, "NTIRE 2017 challenge on single image super-resolution: Dataset and study," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, July 2017.
- [45] D. Ghadiyaram and A. C. Bovik, "Massive online crowdsourced study of subjective and objective picture quality," *IEEE Transactions on Image Processing*, vol. 25, no. 1, pp. 372–387, 2015.
- [46] T. Virtanen, M. Nuutinen, M. Vaahteranoksa, P. Oittinen, and J. 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.
- [47] V. Hosu, H. Lin, T. Sziranyi, and D. Saupe, "Koniq-10k: An ecologically valid database for deep learning of blind image quality assessment," *IEEE Transactions on Image Processing*, vol. 29, pp. 4041–4056, 2020.
- [48] Y. Fang, H. Zhu, Y. Zeng, K. Ma, and Z. Wang, "Perceptual quality assessment of smartphone photography," in *Proceedings of the IEEE/CVF*

- Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3677–3686.
- [49] H. R. Sheikh, M. F. Sabir, and A. C. Bovik, “A statistical evaluation of recent full reference image quality assessment algorithms,” *IEEE Transactions on Image Processing*, vol. 15, no. 11, pp. 3440–3451, 2006.
- [50] 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, pp. 011006–011006, 2010.
- [51] N. Ponomarenko, L. Jin, O. Ieremeiev, V. Lukin, K. Egiazarian, J. Astola, B. Vozel, K. Chehdi, M. Carli, F. Battisti *et al.*, “Image database tid2013: Peculiarities, results and perspectives,” *Signal Processing: Image Communication*, vol. 30, pp. 57–77, 2015.
- [52] H. Lin, V. Hosu, and D. Saupe, “Kadid-10k: A large-scale artificially distorted iqa database,” in *Proceedings of the IEEE Eleventh International Conference on Quality of Multimedia Experience*, 2019, pp. 1–3.
- [53] D. Jayaraman, A. Mittal, A. K. Moorthy, and A. C. Bovik, “Objective quality assessment of multiply distorted images,” in *Proceedings of the Conference Record of the Forty Sixth Asilomar Conference on Signals, Systems and Computers*, 2012, pp. 1693–1697.
- [54] S. Corchs, F. Gasparini, and R. Schettini, “Noisy images-jpeg compressed: Subjective and objective image quality evaluation,” in *Proceedings of the SPIE Image Quality and System Performance XI*, vol. 9016, 2014, pp. 274–282.
- [55] Z. Ni, L. Ma, H. Zeng, J. Chen, C. Cai, and K.-K. Ma, “Esim: Edge similarity for screen content image quality assessment,” *IEEE Transactions on Image Processing*, vol. 26, no. 10, pp. 4818–4831, 2017.
- [56] Z. Ni, H. Zeng, L. Ma, J. Hou, J. Chen, and K.-K. Ma, “A gabor feature-based quality assessment model for the screen content images,” *IEEE Transactions on Image Processing*, vol. 27, no. 9, pp. 4516–4528, 2018.
- [57] Q. Jiang, Y. Gu, C. Li, R. Cong, and F. Shao, “Underwater image enhancement quality evaluation: Benchmark dataset and objective metric,” *IEEE Transactions on Circuits and Systems for Video Technology*, pp. 1–1, 2022.
- [58] N. Yang, Q. Zhong, K. Li, R. Cong, Y. Zhao, and S. Kwong, “A reference-free underwater image quality assessment metric in frequency domain,” *Signal Processing: Image Communication*, vol. 94, p. 116218, 2021.
- [59] C. Li, Z. Zhang, H. Wu, W. Sun, X. Min, X. Liu, G. Zhai, and W. Lin, “Agiqa-3k: An open database for ai-generated image quality assessment,” *arXiv preprint arXiv:2306.04717*, 2023.
- [60] J. Wang, H. Duan, J. Liu, S. Chen, X. Min, and G. Zhai, “Aigciqa2023: A large-scale image quality assessment database for ai generated images: From the perspectives of quality, authenticity and correspondence,” *arXiv preprint arXiv:2307.00211*, 2023.
- [61] W. Yu, X. Zhang, Y. Zhang, Z. Zhang, and J. Zhou, “Blind image quality assessment for a single image from text-to-image synthesis,” *IEEE Access*, vol. 9, pp. 94 656–94 667, 2021.
- [62] S. Su, H. Lin, V. Hosu, O. Wiedemann, J. Sun, Y. Zhu, H. Liu, Y. Zhang, and D. Saupe, “Going the extra mile in face image quality assessment: A novel database and model,” *IEEE Transactions on Multimedia*, 2023.
- [63] G. Hou, Y. Li, H. Yang, K. Li, and Z. Pan, “Uid2021: An underwater image dataset for evaluation of no-reference quality assessment metrics,” *ACM Transactions on Multimedia Computing, Communications and Applications*, vol. 19, no. 4, pp. 1–24, 2023.
- [64] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, “The caltech-ucsd birds-200-2011 dataset,” 2011.
- [65] H. Huang, R. He, Z. Sun, T. Tan *et al.*, “Introvae: Introspective variational autoencoders for photographic image synthesis,” *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [66] Z. Ying, H. Niu, P. Gupta, D. Mahajan, D. Ghadiyaram, and A. Bovik, “From patches to pictures (paq-2-piq): Mapping the perceptual space of picture quality,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3575–3585.
- [67] Z. Pan, F. Yuan, J. Lei, Y. Fang, X. Shao, and S. Kwong, “Vcrnet: Visual compensation restoration network for no-reference image quality assessment,” *IEEE Transactions on Image Processing*, vol. 31, pp. 1613–1627, 2022.
- [68] J. Ke, Q. Wang, Y. Wang, P. Milanfar, and F. Yang, “Musiq: Multi-scale image quality transformer,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5148–5157.
- [69] C. Liu, B. Zoph, M. Neumann, J. Shlens, W. Hua, L.-J. Li, L. Fei-Fei, A. Yuille, J. Huang, and K. Murphy, “Progressive neural architecture search,” in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 19–34.
- [70] A. Gretton, K. Borgwardt, M. Rasch, B. Schölkopf, and A. Smola, “A kernel method for the two-sample-problem,” *Advances in neural information processing systems*, vol. 19, 2006.
- [71] Y. Rubner, C. Tomasi, and L. J. Guibas, “A metric for distributions with applications to image databases,” in *Sixth international conference on computer vision (IEEE Cat. No. 98CH36271)*. IEEE, 1998, pp. 59–66.
- [72] S. Kolouri, K. Nadjahi, U. Simsekli, R. Badeau, and G. Rohde, “Generalized sliced wasserstein distances,” *Advances in neural information processing systems*, vol. 32, 2019.
- [73] S. Kullback and R. A. Leibler, “On information and sufficiency,” *The annals of mathematical statistics*, vol. 22, no. 1, pp. 79–86, 1951.