

# Multi-task Guided Blind Omnidirectional Image Quality Assessment with Feature Interaction

Yun Liu<sup>1</sup>, Sifan Li<sup>1</sup>, Huiyu Duan<sup>1</sup>, Yu Zhou<sup>1</sup>, Daoxin Fan<sup>1</sup>, and Guangtao Zhai<sup>1</sup>

<sup>1</sup>Affiliation not available

July 22, 2024

## Abstract

With the development of virtual reality (VR) applications, omnidirectional image quality assessment (OIQA) has become an increasingly vital problem. In this paper, a multi-task guided blind omnidirectional image quality assessment with local and global feature interaction and fusion is proposed. Specifically, a bidirectional pseudo-reference (BPR) module capturing the error maps on viewports using the two opposite pseudo-reference information is first constructed, which is followed by a multi-scale feature extraction module to obtain multi-scale local degradation features. Moreover, to well complement the local features on viewports, a Mamba module is adopted to extract the multiscale global features. Then the features from the local and global branches are deeply fused based on a multi-level aggregation module. Finally, motivated by the multi-task managing mechanism of human brain, a multi-task learning module is introduced to assist the main quality assessment task. Extensive experimental results demonstrate that our proposed method achieves the state-of-the-art performance on the blind OIQA task compared to other models.

# Multi-task Guided Blind Omnidirectional Image Quality Assessment with Feature Interaction

Yun Liu<sup>1</sup>, Sifan Li<sup>1</sup>, Huiyu Duan<sup>1</sup>, Yu Zhou<sup>1</sup>, Daoxin Fan<sup>1</sup>, and Guangtao Zhai<sup>1</sup>, *Senior Member, IEEE*

**Abstract**—With the development of virtual reality (VR) applications, omnidirectional image quality assessment (OIQA) has become an increasingly vital problem. In this paper, a multi-task guided blind omnidirectional image quality assessment with local and global feature interaction and fusion is proposed. Specifically, a bidirectional pseudo-reference (BPR) module capturing the error maps on viewports using the two opposite pseudo-reference information is first constructed, which is followed by a multi-scale feature extraction module to obtain multi-scale local degradation features. Moreover, to well complement the local features on viewports, a Mamba module is adopted to extract the multi-scale global features. Then the features from the local and global branches are deeply fused based on a multi-level aggregation module. Finally, motivated by the multi-task managing mechanism of human brain, a multi-task learning module is introduced to assist the main quality assessment task. Extensive experimental results demonstrate that our proposed method achieves the state-of-the-art performance on the blind OIQA task compared to other models.

**Index Terms**—Bidirectional pseudo reference, omnidirectional image quality assessment, Mamba, multi-level aggregation, multi-task learning, no-reference (NR).

## I. INTRODUCTION

QUALITY degradation exists in various multimedia contents, which can degrade their perceptual quality and limit the applications. The quality degradation problem is more serious in omnidirectional scenarios due to the increased data volume, and may affect the quality of experience (QoE) more due to the immersive nature of virtual reality (VR) [1]. Therefore, it is significant to develop more effective quality assessment method to further help optimize the QoE in VR environment [2].

As a fundamental problem of QoE assessment, the task of image quality assessment (IQA) has attracted remarkable attention for a long time [3], [4], including two-dimensional (2D) image, omnidirectional image (OI) and so on [5]–[7]. Among them, many omnidirectional image quality assessment

(OIQA) methods have been proposed in recent years with the development of VR area, which include full-reference (FR) metrics [8]–[11], reduced-reference (RR) metrics [12], [13] and no-reference (NR) metrics [14]–[17] according to whether the reference information is introduced. In the early stage, many FR OIQA methods have extended previous 2D IQA models such as peak-signal-to-noise ratio (PSNR) and structural similarity (SSIM) into omnidirectional projection to perform evaluation [8], [18]. Since an omnidirectional image (OI) has a wider viewing range and more complicated perceptual features than a traditional 2D image, the above 2D IQA-based models lead to mediocre performances. In addition, for real applications, it is intractable to obtain the reference image for an omnidirectional image, which makes the NR OIQA metrics have more practical significance than the other two types of OIQA methods.

Motivated by the important role of human visual system (HVS) in image quality assessment area, many NR OIQA models have been proposed by capturing representative visual features and semantic information [19]–[21], which can be categorized into three types, i.e., image-based, patch-based, and viewport-based. Image-based models generally treat an omnidirectional image as a 2D image to capture global visual information. However, due to the wide scene range and large image size, such methods may ignore local visual distortions, which limits the performance [22], [23]. Because the HVS is extremely sensitive to local information, many patch-based models have been proposed based on the cropped 2D image patches from an omnidirectional image, which obtain rich local information to improve the overall performance [24]. Motivated by the viewing characteristics, some viewport-based models have achieved better performance based on the features extracted from each field of view (FoV). Moreover, considering the vital role of local and global features in IQA tasks, some OIQA models have been proposed by integrating local and global information [20], which presents a more effective way to obtain the representative visual features.

As the boost of deep learning, the deep neural network (DNN) based methods have started to show their ability and have become the mainstream of OIQA models [25], [26]. Inspired by the human visual system, Jiang *et al.* [27] have built an effective network by mimicking human visual perception. Xu *et al.* [20] have proposed a GCN based OIQA method with an elaborated viewports choosing algorithm, which also proves the importance of local features. Considering the important role of multi-level features in IQA area, Liu *et al.* [28] have utilized the multi-scale video features to conduct the video quality assessment. Then Sun *et al.*

This work was supported in part by Shenyang science and technology plan project under Grant 23-407-3-32, in part by Liaoning Province Natural Science Foundation under Grant 2023-MS-139, and in part by National Natural Science Foundation of China under Grant 61901205.

Yun Liu, Sifan Li, and Daoxin Fan are with the Faculty of Information, Liaoning University, Shenyang 110036, China (e-mail: yunliu@lnu.edu.cn; sfljohn@foxmail.com; fdx\_0729@163.com).

Huiyu Duan and Guangtao Zhai are with the Institute of Image Communication and Information Processing, Shanghai Key Laboratory of Digital Media Processing and Transmissions, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: huiyuduan@sjtu.edu.cn; zhaiguangtao@sjtu.edu.cn).

Yu Zhou is with the School of Information and Control Engineering, China University of Mining and Technology, Xuzhou 221116, China, and also with Xuzhou First People's Hospital, Xuzhou 221116, China (e-mail: zhouy@cumt.edu.cn).

[29] have proposed a model based on multi-channel CNNs to tackle OIQA challenges and have further improved the overall performance of OIQA. However, few works pay attention to multi-scale local and global information, and their interactive relationship. Although further progress has been achieved by the aforementioned DNN-based OIQA models, significant efforts are still needed to build more effective models.

To deal with the problems and challenges mentioned above, we propose an NR OIQA network by deeply fusing representative multi-scale local and global semantic information. Min *et al.* [30] have introduced a module to generate pseudo reference images to capture the representative distortion information in distorted image, which provides a new way to capture the quality degradation information without the help of reference image. Specifically, inspired by the pseudo reference conception proposed in [30], we first propose a bidirectional pseudo reference module to extract multi-scale local semantic information from the two directions including a restoration direction and a further degradation direction, which can capture important quality changing information to augment the prediction accuracy. Moreover, inspired by Mamba's unprecedented efficiency and capability in extracting features from long range images or texts [31], a global feature extraction module built based on Mamba is developed to obtain multi-scale global information. Then motivated by human hierarchical visual perception characteristics, a multi-level aggregation module is adopted to extract the interactive information and refine the shared features, which can achieve a better performance than the simple fusion or concatenation way. To further optimize the learning process, we further apply a multi-task module to assist the model to adaptively assign weights among different tasks, which can yield a more stable performance. Our main contributions are summarized as follows:

- 1) A bidirectional pseudo reference module is proposed to extract representative local semantic differences from two opposite directions, which can well reflect the local quality degradation and improve the feature representation.
- 2) A Mamba-based feature extraction module is designed to catch efficient multi-scale global visual information to well complement the local features, which can reduce the data volume burden of the model and improve the overall performance.
- 3) A multi-scale interactive feature fusion module is introduced to the OIQA task to strengthen the feature interaction and deep fusion, which can improve the accuracy of our model.
- 4) A multi-task learning module is designed to guide the model to adaptively assign weights among different degradations, which further improves the efficiency of our model.

To illustrate the idea more structurally, we arrange the remainder of this paper as follows. The related works of this paper are briefly reviewed in section II. Our proposed method is introduced in detail in section III, and the experimental results and the analysis are reported in section IV. Finally,

the conclusion of this work is presented in section V.

## II. RELATED WORK

We review the related works including the previous OIQA models and the Mamba structure in this section.

### A. OIQA Models

FR OIQA metrics require all information of the reference image, which can easily obtain the quality difference between a distorted image and a reference image. Many traditional FR OIQA methods extended the previous 2D IQA models to evaluate the quality.

The spherical domain-based model (S-PSNR) was proposed by calculating the PSNR value in the spherical domain. Then a CPP-PSNR metric was designed by calculating the PSNR value in the space of Craster parabolic projection (CPP) [22]. Motivated by the above works, the weighted-to-spherically-uniform PSNR and spherical domain-based SSIM models were built [32], [33]. Although the above FR OIQA models present relatively satisfactory results at the early stage, they are designed based on 2D IQA metrics, and fail to obtain specific visual features of OI, which limits further development of attentive deep image quality assessment for omnidirectional stitching.

FR OIQA model presents a way to capture the quality degradation in distorted image, but the reference information is generally missing in real applications, which makes the NR OIQA models more practical and popular. The existing NR OIQA methods can be categorized into three types, which are the whole image-based methods, patch-based methods and viewport-based methods [19]. The whole image-based methods generally use the equirectangular projection (ERP) image as the input and directly calculate the image quality [23]. The patch-based methods also called as projection space-based methods mainly focus on seeking for a better representation space to obtain more effective features, which digs deeply into the characteristics of the projection methods including the segmented spherical projection (SSP), cube map projection (CMP), equirectangular projection (ERP) [19], etc. For example, in [6], an NR OIQA metric were designed for OIs on the SSP space based on the local details and global features in both bipolar regions of the reprojection space. Jiang *et al.* [27] focused on the local visual features in the CMP space and built a NR OIQA work. Kim *et al.* [34] introduced a blind patch based OIQA metric by segmenting the OI in ERP format into non-overlapping patches in a uniformed size and exploring the positional features of them. Liu *et al.* [35] introduced an effective quality assessment metric by fusing the local structural features and global natural features based on the OI in ERP format.

The viewport-based NR OIQA metrics aim to simulate the visual mechanism when human watching the VR contents. In [20], the final quality score was obtained by calculating both the global prediction quality of the entire image and the local prediction quality of the viewports. Later, considering the importance of viewports in six directions, a multi-channel viewport-based method was designed [29]. Considering the

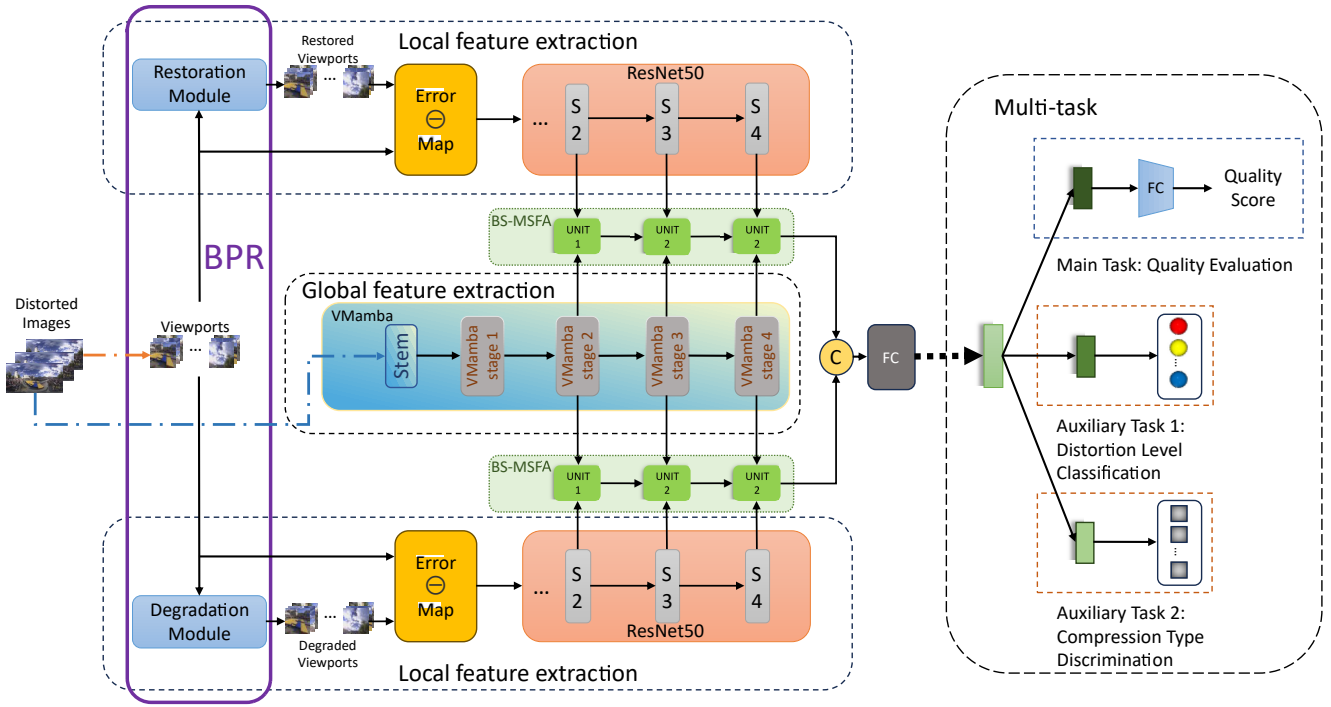


Fig. 1. Overview architecture of our proposed method.

quality degradation in an image is related to the type and degree of the distortion, the auxiliary task for distortion type discrimination was utilized in a multi-stream network [36], which motivated us to build an effective multi-task model. Overall, the above works proves that multi-scale features from both the local and global degrees are particularly important in the OIQA task.

Although the NR OIQA model performs better than the FR and NR models, there is still room for improvement to make the OIQA more effective. More accurate local and global visual features need to be obtained, and the interactive relationship between them should be considered in OIQA models. In addition, it is hard to capture quality changing information without the reference image. An effective module that can capture important quality degradation information hidden in the distorted image should be deeply dug. Overall, OIQA is a complex and challenging work, which needs to take many factors into consideration, such as human visual characteristics, quality degradation between the distorted image and the reference, etc.

### B. Mamba

Mamba [37] has recently drawn considerable attention in various areas, which yields significant results in long sequence modeling tasks [38]–[40]. Mamba consists of repeated Mamba blocks with state space model (SSM) blocks [41]–[43], standard normalization layers, and residual connections, which relieves the constraints of modeling in a convolutional neural network (CNN) [44]–[46]. Compared to Transformer [47], Mamba provides us with advanced and excellent model-

ing capabilities but without secondary computing complexity. Considering the significant advantages over CNNs and Transformer, Mamba demonstrates its enormous potential as a base model for vision tasks, which promotes its further development.

Inspired by Mamba, VMamba [48] was proposed as an efficient model based on down-sampling operations [49] and Visual State Space (VSS) blocks with 2D-selective-scan (SS2D) blocks [31], as shown in Fig. 5. VMamba has bidirectional selective state space model (SSM) blocks [37] along 2D axes by integrating the information from all the other four pixels in different directions around each pixel [31], which can capture rich global semantic information by combining the information of each pixel and reduce the time complexity. Many works have introduced VMamba to various visual tasks and achieved significant performance [50]–[54]. For example, Xie *et al.* [55] introduced VMamba and reformed it in dynamic feature enhancement for multi-modal image fusion, which performs well in medical imaging. Yang *et al.* [56] presented the advantages of Mamba in feature extraction by inventing a scheme with it for image segmentation, which achieved good results. Shi *et al.* [52] introduced a network with VMamba for image restoration and a state-of-the-art performance was achieved. Ma *et al.* [54] then applied VMamba on crowd counting work to solve the problems in counting specific points of a scene, and invented a new approach that inherited the merits of VMamba for global modeling and low computational costs, which achieved a remarkable performance. Considering the remarkable effectiveness of VMamba in the 2D image area, especially its significance for long-range modeling, we choose

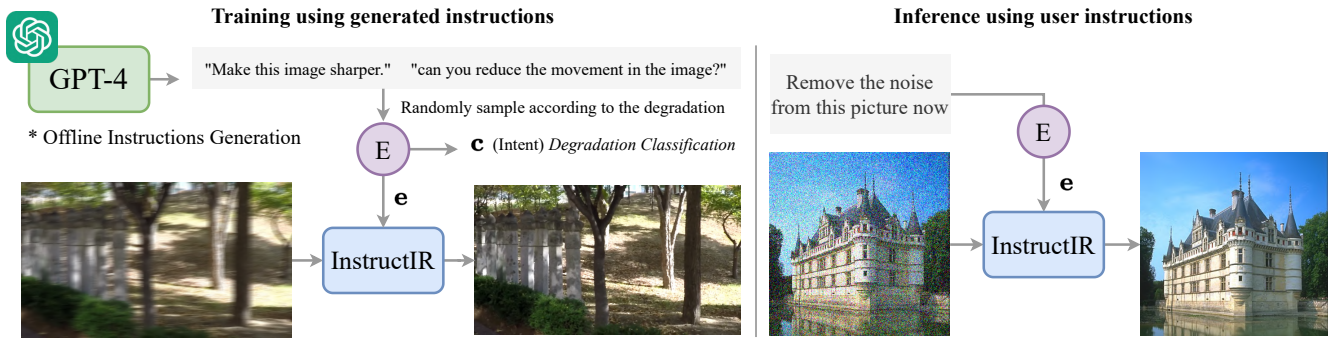


Fig. 2. The method to textually instruct the model to restore images (image from [60]).

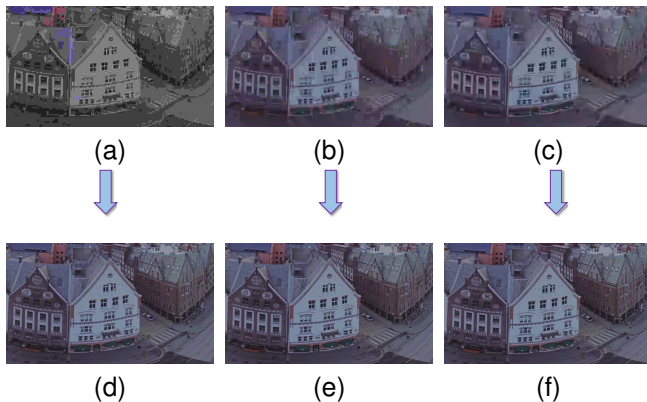


Fig. 3. The distorted viewpoints and their restored viewpoints. (a), (b) and (c) are the viewpoints of the OIs with different types of distortions. (d), (e) and (f) are the restored viewpoints of (a), (b) and (c), respectively.

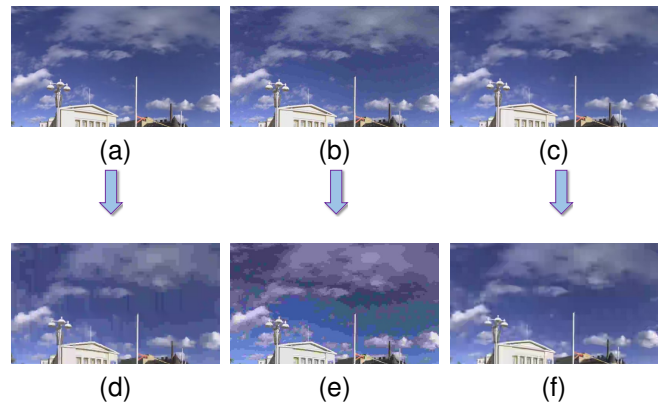


Fig. 4. The distorted viewpoints and their degraded viewpoints. (a), (b) and (c) are the viewpoints of the distorted OIs with different types of distortions (JPEG [68], H.264/AVC [69] and H.265/HEVC [70], respectively). (d), (e) and (f) are the degraded viewpoints of (a), (b) and (c), adding JPEG compression, camera sensor noise and Gaussian noise with  $\sigma = 8$ , respectively.

it as the backbone of the global branch of network to capture the global semantic features.

### III. PROPOSED METHOD

In this section, our proposed model is described in detail. The overall framework is shown in Fig. 1. Specifically, we first extract the representative multi-scale local and global features, respectively. Then an interactive fusion module is built to strengthen the interaction and effective fusion. Finally, a multi-task guided learning module is designed to refine the semantic information that damages image quality and guide the model to adaptively assign weights among different tasks, which can improve our model’s efficiency and performance.

#### A. The Bidirectional Pseudo-Reference (BPR) Module

Motivated by the previous works [19], [57]–[59], we find that the difference between the pseudo-reference (PR) image and the distorted image can reflect the image quality degradation. We introduce a bidirectional pseudo-reference (BPR) module to obtain two pseudo-reference images by restoring and degrading the distorted image, which is presented as the restoration and degradation modules respectively in Fig. 1. For the restoration module as shown in Fig. 2, we adopt a novel generative model [60], which is designed based on GPT-4 and text prompts [61] and achieves an excellent performance

in image restoration. We utilize this model with the textual instructions to generate the restored image by describing *Remove the distortion in the image* or *Restore the quality of the image*. Fig. 3 presents the viewpoints of the distorted OI with different degrees of distortions and their corresponding restored viewpoints. It can be seen that the generated restored viewport has a better quality than the distorted viewport, which proves that the restored viewport can well reflect the quality difference.

For the degradation module, we obtain the degraded images by adding shuffled blur, down-sampling, and noise following the work [62]. Particularly, the down-sampling was randomly used from the nearest, bilinear or bicubic interpolations, and the noise was synthesized by adding Gaussian noise in different levels, JPEG compression or camera sensor noise, which generates the degraded images with random types and random levels of distortions to describe the distorted images in reality. It needs to be mentioned that the overall performance of our model is similar by training the degradation module with a specific type of distortion to random distortions. Since the type and level of the distortion is random in real applications, we train the degradation module by adding random types and random levels of distortions and finally test the model also with the random degradations. Fig. 4 presents the viewpoints

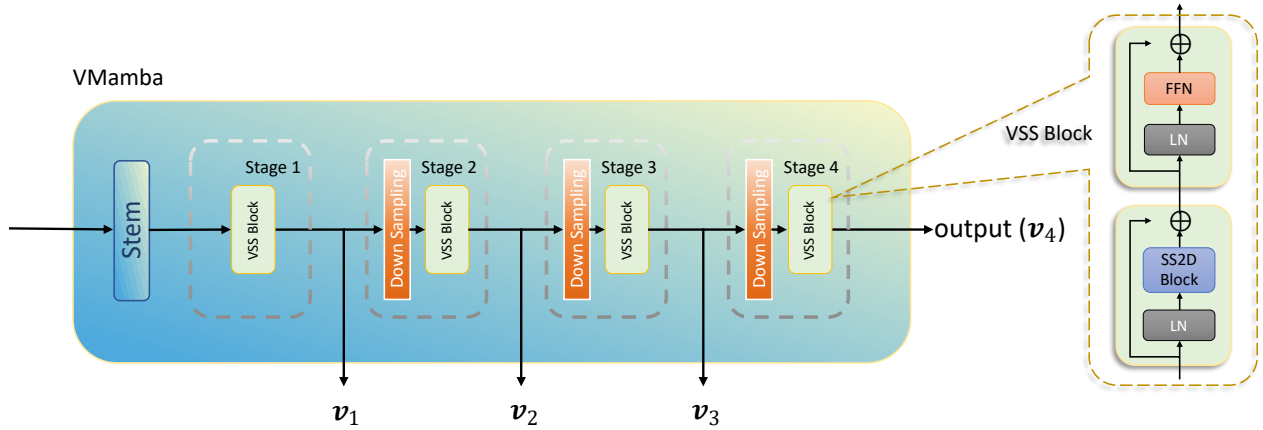


Fig. 5. The structure of VMamba [48] and VSS block.

of the distorted OI with different degrees of distortions and their corresponding degraded viewpoints. Three types of distortions with three different levels are randomly aggravated to the inputs to generate the degraded images that has worse qualities. Other than restored images that present the reference score from a positive degree, degraded PR images can help the prediction model compare with the worst quality score, which can reflect the quality difference.

### B. Local Feature Extraction

Considering that only one local viewport of an omnidirectional image is watched at a time for a specific user, 20 viewports are first generated following the work in [57] for the local feature extraction. Then these 20 viewports are fed into the BPR module to obtain the pseudo-reference (PR) viewports from two opposite directions, i.e., a restoration direction and a degradation direction. To effectively capture the presentative spatial quality changing between the PR viewports and the distorted viewports, error maps are calculated to capture the rich semantic difference information from the two directions. Here, in order to obtain a better correlation with the quality perceived by viewers, the normalized log difference function [58] is adopted, which is defined as in (1):

$$E = \log_{\alpha} (\alpha + (I_r - I_d)^2), \quad (1)$$

where  $\alpha = \epsilon/255^2$  is a constant with  $\epsilon = 0.1$ ,  $I_r$  is the value of each pixel of the PR viewport, and  $I_d$  is the value of each pixel of the distorted viewport. After the necessary procedure, the error map is then fed into the local feature extraction module. In this paper, we elaborately take ResNet50 as the backbone to obtain the local features. The ResNet50 consists of five stages: stage 0, 1, 2, 3 and 4, respectively, in which stage 1 includes three bottlenecks, and stage 2, 3 and 4 include 4, 6 and 3 bottlenecks, respectively [63]. We take the features from the last three stages as local features in our method.

### C. Global Feature Extraction

Since VMamba recently presents a remarkable performance in feature extraction in deep learning [48], we take it as the

backbone of the global feature extraction module, shown in Fig. 5. For an OI, it is firstly partitioned into patches with a stem module, and a 2D feature map  $M_d \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4}}$  with the spatial dimension of  $\frac{H}{4} \times \frac{W}{4}$  [48] where  $H$  and  $W$  are the height and width of the image is consequently obtained to feed into the VMamba. Denoting the output of stage  $n$  as  $v_n$ , the output of the first stage can be obtained as in (2) and (3):

$$T_1 = \text{SS2DB}(\text{Norm}(M_d)) \oplus M_d, \quad (2)$$

$$v_1 = \text{FFN}(\text{LN}(T_1)) \oplus T_1. \quad (3)$$

For the output of stage  $n$  ( $n > 1$ ) can be obtained as in (4) and (5):

$$T_n = \text{SS2DB}(\text{Norm}(\text{DS}(v_{n-1}))) \oplus \text{DS}(v_{n-1}), \quad (4)$$

$$v_n = \text{FFN}(\text{LN}(T_n)) \oplus T_n, \quad (5)$$

where  $\text{SS2DB}(\cdot)$  denotes the SS2D block [48],  $\text{Norm}(\cdot)$  denotes the normalization layer,  $\text{FFN}(\cdot)$  is a feedforward neural network,  $\text{LN}(\cdot)$  means layer normalization, and  $\text{DS}(\cdot)$  is for down-sampling operation. The SS2D block is defined as in (6) and (7):

$$F_{\text{SS2D}}(\cdot) = \text{SS2D}(\text{SiLU}(\text{Dwc}(\text{Linear}(\cdot))))), \quad (6)$$

$$\text{SS2DB}(\cdot) = \text{Linear}(\text{LN}(F_{\text{SS2D}}(\cdot))), \quad (7)$$

where  $\text{SS2D}(\cdot)$  is the 2D-Selective-Scan operation [64], and  $\text{Dwc}(\cdot)$  is a  $3 \times 3$  depth-wise convolution layer. The outputs  $v_2$ ,  $v_3$ , and  $v_4$  are then prepared for the multi-scale fusion.

### D. Interactive Fusion Module

To fully apply the interactive relationship between local and global features, we adopt an interactive fusion module to extract the interactive information and refine the shared features. Due to the task insensitivity, the regression learning module based on the output of one layer fails to perform a satisfactory performance, while the multi-scale representation provides us

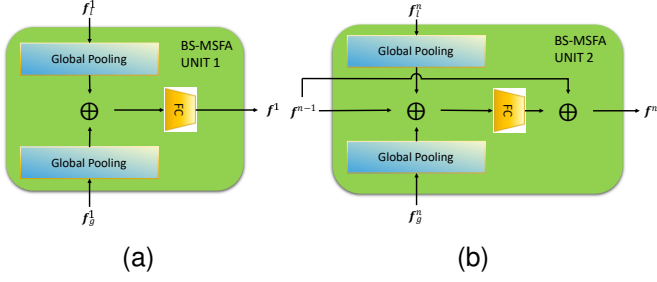


Fig. 6. BS-MSFA module: (a) is BS-MSFA UNIT 1, and (b) is BS-MSFA UNIT 2.

with a new way to capture the crucial features [28], [65]. Here, we apply the outputs of last three stages of ResNet50 as the multi-scale local features and take the outputs of the last three Mamba stages as the multi-scale global features. To avoid the rigid connection between multi-level features, the interactive fusion module with residual structure, namely Bi-Stream Multi-Scale Feature Aggregation (BS-MSFA) module, is proposed to reduce the shallow feature dimension and achieve efficient fusion by digging the interactive relationship between local and global features, shown in Fig. 6.

Specifically, BS-MSFA consists of a BS-MSFA unit 1 and two BS-MSFA unit 2's. For the unit 1, two features from both branches are fed into it to capture the first level fusion features  $f^1$ , which is defined as in (8):

$$f^1 = \text{FC}(g(f_l^1) \oplus g(f_g^1)), \quad (8)$$

where  $\text{FC}(\cdot)$  denotes a fully-connected layer,  $g(\cdot)$  denotes a global pooling operation, and  $\oplus$  means the concatenate operation.  $f_l^1$  and  $f_g^1$  are the first level of local features and global features, respectively.

Then the first level fusion features  $f^1$ , the second level local features and the second level global features are then fed into the BS-MSFA unit 2 to capture the second level fusion features  $f^2$ . Before feeding  $f^1$  into the BS-MSFA unit 2,  $f^1$  is processed with a fully-connected layer. The second level fusion feature  $f^2$  is shown as in (9):

$$f^2 = f^1 \oplus \text{FC}(f^1 \oplus g(f_l^2) \oplus g(f_g^2)). \quad (9)$$

The same as the above procedure, the third level fusion features are then obtained. Then the multi-scale fusion features of the 20 viewports from the restoration viewports and degradation viewports are concatenated together and then processed by a fully-connected layer for the multi-task module.

### E. The Multi-Task Module

Considering that the compression types and distortion degrees have different impacts on quality perception, we design a multi-task module to refine the shared information and optimize the model's performance by the assistance of the distortion level classification task and compression type discrimination task. For the main task, a fully-connected layer is adopted to get the final quality score. For the two auxiliary tasks, two fully connected layers containing 1024 nodes and 64 nodes are utilized.

### F. Network Training

For our multi-task method, Euclidean loss function and cross-entropy loss function are adopted to optimize the network [36]. The Euclidean loss is defined as in (10):

$$L_q(V_k; W_0) = \frac{1}{N} \sum_{k=1}^N \|s_k(V_k) - q_k(W_0)\|_2^2, \quad (10)$$

where  $k$  denotes the ordinal number of the  $k$ -th training sample.  $s_k$  is the subjective quality score and  $q_k$  is the predicted quality score.  $V_k$  means the viewport, and  $W_0$  denotes the parameters.

The cross-entropy loss based on the distortion level classification task is defined as in (11):

$$L_d(V_k; W_0) = - \sum_{k=1}^N \sum_{i=1}^M m_k^i \log \hat{p}_k^i(V_k; W_0), \quad (11)$$

where  $m_k^i$  is the ground-truth multi-class indicator vector. If the  $k$ -th training sample is in the  $i$ -th distortion level, then  $m_k^i$  will be one, otherwise  $m_k^i$  will be zero.  $\hat{p}_k^i$  denotes the predicted probability of whether the distortion in the  $k$ -th training sample is in the  $i$ -th distortion level.

Similar to the loss for the distortion level classification, the cross-entropy loss of the compression type discrimination task is defined as in (12):

$$L_c(V_k; W_0) = - \sum_{k=1}^N \sum_{i=1}^C c_k^i \log \hat{r}_k^i(V_k; W_0), \quad (12)$$

where  $c_k^i$  is the ground-truth multi-class indicator vector for the compression type discrimination. If the distortion type of the  $k$ -th training sample is the  $i$ -th compression type, the  $c_k^i$  will be one, otherwise it will be zero.  $\hat{r}_k^i$  denotes the predicted probability of whether the compression type of the  $k$ -th training sample is the  $i$ -th compression type.

To learn the superior parameters for all the three tasks, the total loss is defined as in (13):

$$L = L_q + 0.1L_d + 0.1L_c, \quad (13)$$

where the correspondences of  $L_q$ ,  $L_d$  and  $L_c$  denote the importance of the three losses and are empirically set to 1, 0.1 and 0.1 in this work, respectively.

## IV. EXPERIMENTAL RESULTS

### A. Experimental Settings

1) **Datasets: OIQA database** [66]: It consists of 336 omnidirectional images in equirectangular format, which consists of 320 distorted OIs based on two types of compressions (JPEG and JPEG2000) and two types of degradations (Gaussian blur (GB), and Gaussian noise (GN)), and 16 reference images. The resolutions of the OIs vary in a range from  $11332 \times 5666$  to  $13320 \times 6660$ . The Mean Opinion Score (MOS) value of each image is provided with the dataset by conducting the subjective experiment in which the single-stimulus (SS) method [67] is adopted.

**CVIQ database** [29]: It provides 544 images in total, which include 528 distorted omnidirectional images based on three different types of coding compressions (JPEG [68], H.264/AVC [69] and H.265/HEVC [70]). The rest 16 OIs are the reference images. All the images are in the same resolution of  $4096 \times 2048$ . Like [66], the SS method is also adopted in the subjective experiment to get the MOS value.

**OSIQA database** [71]: It is a database designed for omnidirectional stitching image quality assessment, which is used to validate the generalization ability of our model. It provides 700 omnidirectional stitching images generated by stitching the packs from different views, which includes 350 distorted based on 14 scenes with different stitching distortions. MOS values are provided on the basis of extensive experiments in [71].

2) *Evaluation Criteria*: Three prevalent criteria which are Pearson Linear Correlation Coefficient (PLCC), Spearman Rank-order Correlation Coefficient (SRCC) and Root Mean Squared Error (RMSE) are adopted to make the monotonicity comparison and accuracy prediction. The three criteria are formulated as in (14), (15) and (16):

$$\text{PLCC} = \frac{\sum_{i=1}^N (s_i - \bar{s})(p_i - \bar{p})}{\sqrt{\sum_{i=1}^N (s_i - \bar{s})^2 \sum_{i=1}^N (p_i - \bar{p})^2}}, \quad (14)$$

where  $N$  denotes the number of the samples.  $s_i$  is the MOS of the  $i$ -th sample, and  $p_i$  is the prediction score.  $\bar{s}$  is the mean value of the MOS's, and  $\bar{p}$  is the mean value of the score that the model predicted for each sample.

$$\text{SRCC} = 1 - \frac{6 \sum_{i=1}^N d_i^2}{N(N^2 - 1)}, \quad (15)$$

where  $d_i$  denotes the distance between the rank of the MOS and the rank of the prediction score given by the model for the  $i$ -th sample.

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (s_i - p_i)^2}, \quad (16)$$

where  $s_i$  is the MOS of the  $i$ -th sample and  $p_i$  is the prediction score given by the model.

Among the three criteria, PLCC and RMSE are calculated by the five-parameter nonlinear mapping [72], [73], which aims to unify the prediction scores given by different metrics into the same range [36].

3) *Comparison of the Metrics*: To validate the advantages of our proposed method, many state-of-the-art methods, including FR and NR models, are employed for the comparison. Many models that were originally designed to take expertise for omnidirectional image quality assessment are also employed for the comparison. Particularly, the FR metrics including PSNR, SSIM [8], S-PSNR [18], CPP-PSNR [22] and WS-PSNR [32] are used. The NR metrics including BRISQUE [74], DESQUE [23], dipIQ [25], MEON [26], BMPRI [30], SSP-BOIQA [6], MC360IQA [29], Zhou *et al.* [36], VGCN [20] and PICS (Pro.) [19], among which, S-PSNR [18], CPP-PSNR [22], WS-PSNR [32], SSP-BOIQA [6], MC360IQA

[29], Zhou *et al.* [36], VGCN [20] and PICS (Pro.) [19] specifically designed for OIQA are employed.

4) *Implementation Details*: In implementation, the dataset is split into a training set and a testing set following the commonly used standard method in [75]–[78]. PyTorch framework [79] is adopted to implement the proposed method and the fine-tuning operation is implemented on both the OIQA and the CVIQ datasets. The SGD optimization [80]–[82] is employed with the momentum parameter set to 0.9, while the batch size and the weight decay parameter are set to 16 and  $10^{-4}$ , respectively. For the two branches of the network, we set the initial learning rates to  $10^{-3}$ , and the learning rate drops with a factor of 0.9 for each epoch with the total number of epochs is 300. The entire experiment processes are implemented on a device with Intel(R) Core (TM) i7-10870H CPU with 16 GB RAM, and one NVIDIA GeForce RTX 2060 graphic card, which shows the advantage of the proposed method on lower-level devices.

## B. Performance Evaluation

1) *Comparison with State-of-the-Arts*: To prove the advanced performance of our model, several OIQA models have been adopted to conduct performance comparisons. The experimental results of all the metrics on the OIQA dataset and CVIQ dataset are summarized in Table I and II, respectively. The top performances are emphasized with boldface. The results demonstrated that the latest deep learning-based NR OIQA models without any reference information, such as MC360IQA [29], SSP-BOIQA [6], Zhou *et al.* [36], VGCN [20], and PICS [19], achieve a better overall performance than all the FR quality metrics and some early NR OIQA models. One of the reasons could be that FR methods mainly rely on handcrafted features, and deep learning technology matures gradually, which boosts the improvement of OIQA models. Although SSIM model [8] presents a promising performance on the whole dataset and on each distortion type, it highly relies on the reference information without considering human visual characteristics on OI, which limited its application.

For the NR metrics in the comparison experiment, it can be found that BMPRI [30] with multiple pseudo reference images (MPRI) takes the worst overall performance among all the NR metrics. One probable reason is the big intervals between two degradation degrees, which may make the network puzzled in building a clear reference standard and lead to a deficient performance. In addition, this model yields the worst results on JP2K and GB distortion among all the models, since it is impossible to design one pseudo generation model to cover all types of distortions, and it is hard to obtain an accurate result only relying on the pseudo quality changing information. So, a more reasonable way to build an effective OIQA model should be capturing the pseudo reference information from different directions and training the network with them and notable features, which are one of our contributions. Compared to MC360IQA [29], SSP-BOIQA [6], Zhou *et al.* [36], and VGCN [20], PICS [19] achieves a better performance, which may benefit from the generative complementary images to fill the semantic blanks of an OI. However, it ignored the

TABLE I  
THE PERFORMANCE COMPARISON ON THE OIQA DATABASE [66]

Metrics	JPEG			JP2K			GN			GB			Overall			
	PLCC	SRCC	RMSE	PLCC	SRCC	RMSE	PLCC	SRCC	RMSE	PLCC	SRCC	RMSE	PLCC	SRCC	RMSE	
FR	PSNR	0.758	0.731	10.245	0.781	0.768	9.379	0.958	0.931	3.654	0.529	0.506	11.268	0.492	0.497	12.528
	SSIM [8]	0.803	0.934	9.355	0.802	0.936	8.985	0.904	0.886	5.467	0.768	0.925	8.500	0.856	0.880	7.436
	S-PSNR [18]	0.87	0.829	7.738	0.816	0.849	8.686	0.919	0.885	5.033	0.699	0.692	9.501	0.716	0.712	10.030
	CPP-PSNR [22]	0.865	0.829	7.873	0.849	0.837	7.943	0.920	0.885	5.001	0.672	0.667	9.830	0.707	0.703	10.167
	WS-PSNR [32]	0.861	0.828	7.994	0.844	0.832	8.070	0.922	0.885	4.942	0.661	0.658	9.966	0.689	0.693	10.428
NR	BRISQUE [74]	0.935	0.921	8.689	0.725	0.733	11.355	0.968	0.979	4.551	0.844	0.857	9.161	0.823	0.831	9.262
	DESQUE [23]	0.897	0.868	6.952	0.739	0.732	10.120	0.953	0.937	3.882	0.749	0.663	8.799	0.725	0.712	9.903
	dipIQ [25]	0.829	0.789	8.783	0.916	0.918	6.030	0.955	0.943	3.772	0.932	0.898	4.816	0.701	0.691	10.259
	MEON [26]	0.823	0.779	8.935	0.680	0.601	11.017	0.952	0.930	3.895	0.764	0.716	8.572	0.749	0.717	9.536
	BMPRI [30]	0.918	0.909	6.210	0.185	0.166	14.768	0.961	0.949	3.534	0.356	0.354	12.248	0.431	0.338	12.984
	SSP-BOIQA [6]	0.877	0.834	7.620	0.853	0.852	7.501	0.905	0.843	5.451	0.854	0.862	6.834	0.860	0.865	7.313
	MC360IQA [29]	0.912	0.901	6.535	0.896	0.882	6.573	0.913	0.926	5.240	0.893	0.918	6.072	0.890	0.909	6.697
	Zhou <i>et al.</i> [36]	0.936	0.94	5.691	0.920	0.934	5.886	0.968	0.957	3.330	0.925	0.920	4.972	0.899	0.923	6.396
	VGCN [20]	0.954	0.929	4.288	0.977	0.946	4.313	0.981	0.975	3.617	0.985	0.965	4.213	0.958	0.952	4.385
	PICS (Pro.) [19]	0.968	0.946	3.988	0.980	0.972	4.047	0.989	0.983	3.575	0.990	0.974	3.827	0.970	0.964	3.991
	<b>Ours</b>	<b>0.979</b>	<b>0.971</b>	<b>3.078</b>	<b>0.990</b>	<b>0.982</b>	<b>3.558</b>	<b>0.995</b>	<b>0.987</b>	<b>1.998</b>	<b>0.992</b>	<b>0.985</b>	<b>2.961</b>	<b>0.992</b>	<b>0.982</b>	<b>2.641</b>

TABLE II  
THE PERFORMANCE COMPARISON ON THE CVIQ DATABASE [29]

Metrics	JPEG			H.264/AVC			H.265/HEVC			Overall			
	PLCC	SRCC	RMSE	PLCC	SRCC	RMSE	PLCC	SRCC	RMSE	PLCC	SRCC	RMSE	
FR	PSNR	0.889	0.766	7.824	0.784	0.783	7.674	0.746	0.745	8.000	0.786	0.757	8.692
	SSIM [8]	0.852	0.929	8.946	0.941	0.940	4.177	0.918	0.917	4.763	0.897	0.885	6.230
	S-PSNR [18]	0.892	0.778	7.727	0.789	0.786	7.589	0.762	0.758	7.785	0.785	0.761	8.714
	CPP-PSNR [22]	0.884	0.765	7.996	0.779	0.777	7.751	0.751	0.748	7.936	0.779	0.754	8.822
	WS-PSNR [32]	0.880	0.756	8.101	0.775	0.773	7.814	0.747	0.744	7.993	0.777	0.751	8.850
NR	BRISQUE [74]	0.913	0.938	5.144	0.780	0.779	7.715	0.771	0.758	8.340	0.826	0.828	7.572
	DESQUE [23]	0.912	0.870	7.003	0.385	0.173	11.410	0.328	0.152	11.362	0.566	0.417	11.603
	dipIQ [25]	0.928	0.793	6.353	0.620	0.635	9.695	0.361	0.326	11.216	0.706	0.623	9.960
	MEON [26]	0.808	0.566	10.057	0.599	0.574	9.900	0.783	0.782	7.484	0.665	0.567	10.510
	BMPRI [30]	0.776	0.498	10.767	0.533	0.520	10.459	0.846	0.840	6.412	0.627	0.621	10.962
	SSP-BOIQA [6]	0.915	0.853	6.847	0.885	0.861	7.042	0.854	0.841	6.302	0.890	0.856	6.941
	MC360IQA [29]	0.941	0.923	5.804	0.932	0.941	5.357	0.914	0.899	4.801	0.939	0.904	4.606
	Zhou <i>et al.</i> [36]	0.957	0.961	5.601	0.953	0.949	3.873	0.929	0.914	4.525	0.902	0.911	6.117
	VGCN [20]	0.989	0.976	2.359	0.972	0.966	3.149	0.940	0.943	4.026	0.965	0.964	3.657
	PICS (Pro.) [19]	0.990	0.983	2.136	0.976	0.972	2.967	0.959	0.962	3.577	0.976	0.973	3.290
	<b>Ours</b>	<b>0.992</b>	<b>0.993</b>	<b>1.953</b>	<b>0.983</b>	<b>0.974</b>	<b>2.556</b>	<b>0.995</b>	<b>0.987</b>	<b>3.053</b>	<b>0.987</b>	<b>0.991</b>	<b>2.734</b>

importance of multi-scale fusion of the features. Our model not only captures the quality changing information from two opposite directions, but also interactively fuses multi-level local and global features, which achieves the best results not only on the whole dataset, but also on each distortion type of the dataset. Overall, the results above prove that our model is reasonable, and able to be effectively applied to evaluate the quality of OIs.

To further prove the practical effectiveness of our model, Table II gives a performance comparison on each compression type and the overall CVIQ database. It can be observed that all the models present a better performance, and the trend of the performances in Table II present a similar trend to those in Table I. Our model also has the best performance on the overall CVIQ dataset and each distortion type. It needs to be mentioned that CVIQ dataset is focused on the types of compression distortion, which is more useful for image transcoding or transmission. Our model achieves the best performance, which further demonstrates the superiority and the potential practical application of our model.

2) *Effect of the Training-Testing Proportion*: To study the effect of the training-testing proportion, experiments on how

the performance varies with the different dataset proportions are conducted. The procedure is repeated five times, and the average results are presented in Table III. It can be seen that as the number of training samples increases, the performances of all three metrics are improved. It needs to be mentioned that the performance of our proposed method is superior to VGCN [20] and PICS [19] methods on all training-testing splits. Moreover, a remarkable performance is still achieved by our method even though only half of the images are employed for training. Specifically, both the PLCC and SRCC are over 0.8, which performs even better than some of the existing metrics with the best proportion in Table I and II. Furthermore, it can be concluded that our method is relatively dependent on the quantity of the images for model training, which means an application potential.

3) *Generalization Ability*: In order for validating the generalization ability and the robustness of our method, the cross-database validation is conducted. We first train the model on one of the two datasets and then test on the other. The results of the cross-database validation are presented in Table IV. It can be observed that all the results are lower than the results of training and testing on the same datasets. The PICS (Pro.)

TABLE III

PERFORMANCE OF OUR PROPOSED METHOD AND TWO ADVANCED METHODS (VGCN [20] AND PICS (PRO.) [19]) WITH DIFFERENT QUANTITIES OF TRAINING SAMPLES ON OIQA DATABASE AND CVIQ DATABASE

Database	Proportion	VGCN			PICS (Pro.)			Ours		
		PLCC	SRCC	RMSE	PLCC	SRCC	RMSE	PLCC	SRCC	RMSE
OIQA	0.8/0.2	0.958	0.952	4.385	0.970	0.964	3.991	0.992	0.982	2.641
	0.7/0.3	0.868	0.870	7.254	0.896	0.893	6.483	0.909	0.910	5.755
	0.6/0.4	0.802	0.793	8.739	0.817	0.806	8.664	0.844	0.842	7.864
	0.5/0.5	0.725	0.721	9.918	0.753	0.737	9.315	0.792	0.779	7.839
	0.4/0.6	0.673	0.652	10.826	0.682	0.654	10.634	0.721	0.699	9.497
CVIQ	0.8/0.2	0.965	0.964	3.657	0.976	0.973	3.290	0.992	0.985	2.961
	0.7/0.3	0.903	0.805	5.941	0.915	0.911	5.307	0.947	0.941	4.909
	0.6/0.4	0.831	0.819	7.796	0.844	0.835	7.685	0.900	0.880	6.908
	0.5/0.5	0.755	0.730	9.381	0.795	0.782	8.241	0.821	0.816	6.918
	0.4/0.6	0.694	0.658	10.024	0.713	0.707	9.846	0.799	0.776	8.572

TABLE IV

THE RESULTS OF THE CROSS-DATABASE VALIDATION

Train	Test	Criterion	BRISQUE	dipIQ	MEON	BMPRI	MC360IQA	Zhou <i>et al.</i>	VGCN	PICS (Pro.)	Ours
CVIQ	OIQA	PLCC	0.682	0.583	0.604	0.331	0.705	0.735	0.787	0.827	<b>0.905</b>
		SRCC	0.524	0.502	0.551	0.192	0.684	0.684	0.778	0.815	<b>0.903</b>
		RMSE	10.870	11.747	11.399	13.576	10.178	10.178	5.437	5.124	<b>4.440</b>
OIQA	CVIQ	PLCC	0.754	0.630	0.688	0.586	0.823	0.823	0.924	0.935	<b>0.970</b>
		SRCC	0.689	0.587	0.624	0.548	0.814	0.814	0.901	0.931	<b>0.972</b>
		RMSE	9.381	10.904	10.145	11.403	7.811	7.811	5.462	4.887	<b>3.857</b>

[19] metric presents a potential competition but is far inferior to our method, although none of the other seven metrics can compete against it on cross-database performance. Our method achieves the highest PLCC and SRCC and the lowest RMSE and is the only method performing with both the PLCC and SRCC over 0.9 when the model is trained on CVIQ dataset and tested on OIQA dataset. In addition, it can be observed that the results of models trained on OIQA dataset are all better than the models trained on CVIQ dataset. One probable reason is that the OIQA database includes four types of compressions and degradations, while the CVIQ database includes images with only the compression types of distortions.

To further prove the generalization ability of our method, the omnidirectional stitching image quality assessment (OSIQA) [71] task is conducted, and the results is shown in Table V. It can be seen that, even compared to the state-of-the-art in the OSIQA area, our method achieves a remarkable performance, which proves the potential performance in the OSIQA area. Our model achieves the best performance on PLCC and SRCC criteria and ranks second on RMSE. The probable reason is that the geometry distortion in OSIQA tasks is not considered, but the RMSE value of our model is close to the top one, which proves the generalization ability of our method.

Overall, the results above indicate the generalization ability of our proposed method.

### C. Ablation Experiments

To verify the contribution of the local feature, the global feature, the bi-stream multi-scale feature aggregation (BS-MSFA) module, the auxiliary tasks module, and replacing the VMamba with Transformer, respectively, the ablation experi-

TABLE V  
PERFORMANCE COMPARISON OF THE STATE-OF-THE-ART NR-IQA MODELS ON OSIQA DATABASE

Database	Model \ Criteria	PLCC	SRCC	RMSE
OSIQA	BRISQUE [74]	0.3072	0.2450	12.296
	NIQE [83]	0.3167	0.2288	12.053
	CORNIA [84]	0.3404	0.2271	12.008
	QAC [85]	0.5747	0.2635	9.9765
	ILNIQE [86]	0.3957	0.1658	11.707
	LPSI [87]	0.5789	0.2127	10.599
	HOSA [88]	0.3270	0.2457	11.859
	dipIQ [25]	0.2394	0.1994	499.01
	BPRI [89]	0.5993	0.2656	9.9980
	BPRI-LSS [89]	0.4889	0.3200	10.994
	BPRI-PSS [89]	0.5085	0.2356	10.957
	BPRic [89]	0.5685	0.3171	10.270
	BMPRI [30]	0.3703	0.2666	11.320
	MC360IQA [29]	0.7943	0.6807	6.9597
	OSIQA-NR [71]	0.8214	0.7236	<b>6.2442</b>
	<b>Ours</b>	<b>0.8670</b>	<b>0.7541</b>	6.8241

ments are conducted, and the results are presented in Table VI. It can be concluded that each component has contribution to the final performance, while the model without the local branch yields the worst performance, which proves that the local branch is the most essential information for OIQA. The model without the Mamba module ranks second to the last, which proves that it is reasonable to combine the local and global features. The model without BS-MSFA module designed by using the traditional concatenating fusion way presents a worse performance than the proposed model, which

TABLE VI  
THE RESULTS OF THE ABLATION EXPERIMENTS

Module	Status	OIQA			CVIQ		
		PLCC	SRCC	RMSE	PLCC	SRCC	RMSE
Local Branch	✓	<b>0.992</b>	<b>0.982</b>	<b>2.641</b>	<b>0.987</b>	<b>0.991</b>	<b>2.734</b>
	✗	0.915	0.921	12.669	0.934	0.923	9.600
VMamba	✓	<b>0.992</b>	<b>0.982</b>	<b>2.641</b>	<b>0.987</b>	<b>0.991</b>	<b>2.734</b>
	✗	0.944	0.948	4.722	0.951	0.954	4.591
BS-MSFA	✓	<b>0.992</b>	<b>0.982</b>	<b>2.641</b>	<b>0.987</b>	<b>0.991</b>	<b>2.734</b>
	✗	0.969	0.975	3.227	0.976	0.972	3.141
Multi-task	✓	<b>0.992</b>	<b>0.982</b>	<b>2.641</b>	<b>0.987</b>	<b>0.991</b>	<b>2.734</b>
	✗	0.982	0.968	3.193	0.983	0.974	3.175
Global Branch	VMamba	<b>0.992</b>	<b>0.982</b>	<b>2.641</b>	<b>0.987</b>	<b>0.991</b>	<b>2.734</b>
Backbone	ViT	0.981	0.979	4.499	0.975	0.974	3.133

presents the significant role of the bi-stream multi-scale feature aggregation. And the model without the assistance of multi-task module ranks the second lowest, which proves the positive effect of two auxiliary tasks. The model using a Transformer network as the backbone to capture the global features presents a worse performance than our final model, which means the effectiveness of Mamba in OIQA tasks. All the above results prove all our contributions and the reasonableness of our model. Overall, our proposed model can be effectively applied to evaluate the OI quality and achieve a high consistency with the human perception.

## V. CONCLUSION

In this paper, a multi-task framework based on multi-scale local and global features is proposed for OIQA. Considering the quality degradation varying in distorted images, a bidirectional pseudo reference module is utilized to capture the rich local features from two opposite directions. Based on the state-of-art performance of Mamba in features extraction, we adopt it as the global feature extractor to obtain multi-scale global information, which can complement the above local features well. To utilize the interactive relationship between the local and global information, the multi-scale feature aggregation module is constructed to make a hierarchically deep fusion. Furthermore, a multi-task learning is applied to optimize the entire model for the quality prediction. The experimental results demonstrate that our method can effectively and accurately predict the quality of an omnidirectional image. In the future, we will extend our method to more omnidirectional quality assessment tasks, and will develop more effective models with advanced technology to evaluate the quality of omnidirectional images by digging into the relationship between human visual characteristics and high-level semantic information.

## REFERENCES

- [1] X. Min, H. Duan, W. Sun, Y. Zhu, and G. Zhai, "Perceptual video quality assessment: A survey," 2024.
- [2] H. Duan, X. Zhu, Y. Zhu, X. Min, and G. Zhai, "A quick review of human perception in immersive media," *IEEE Open Journal on Immersive Displays*, vol. 1, pp. 41–50, 2024.
- [3] Z. Wang, E. Simoncelli, and A. Bovik, "Multiscale structural similarity for image quality assessment," in *The Thirty-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*, vol. 2, 2003, pp. 1398–1402 Vol.2.
- [4] H. Sheikh and A. Bovik, "Image information and visual quality," *IEEE Transactions on Image Processing*, vol. 15, no. 2, pp. 430–444, 2006.
- [5] M. Xu, C. Li, Z. Chen, Z. Wang, and Z. Guan, "Assessing visual quality of omnidirectional videos," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 12, pp. 3516–3530, 2019.
- [6] X. Zheng, G. Jiang, M. Yu, and H. Jiang, "Segmented spherical projection-based blind omnidirectional image quality assessment," *IEEE Access*, vol. 8, pp. 31 647–31 659, 2020.
- [7] G. Yue, C. Hou, T. Zhou, and X. Zhang, "Effective and efficient blind quality evaluator for contrast distorted images," *IEEE Transactions on Instrumentation and Measurement*, vol. 68, no. 8, pp. 2733–2741, 2019.
- [8] Z. Wang, A. Bovik, H. Sheikh, and E. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [9] L. Zhang, Y. Shen, and H. Li, "Vsi: A visual saliency-induced index for perceptual image quality assessment," *IEEE Transactions on Image Processing*, vol. 23, no. 10, pp. 4270–4281, 2014.
- [10] 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.
- [11] X. Sui, K. Ma, Y. Yao, and Y. Fang, "Perceptual quality assessment of omnidirectional images as moving camera videos," *IEEE Transactions on Visualization and Computer Graphics*, vol. 28, no. 8, pp. 3022–3034, 2022.
- [12] M. Narwaria, W. Lin, I. V. McLoughlin, S. Emmanuel, and L.-T. Chia, "Fourier transform-based scalable image quality measure," *IEEE Transactions on Image Processing*, vol. 21, no. 8, pp. 3364–3377, 2012.
- [13] J. Wu, W. Lin, G. Shi, and A. Liu, "Reduced-reference image quality assessment with visual information fidelity," *IEEE Transactions on Multimedia*, vol. 15, no. 7, pp. 1700–1705, 2013.
- [14] Q. Wu, H. Li, K. N. Ngan, and K. Ma, "Blind image quality assessment using local consistency aware retriever and uncertainty aware evaluator," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 9, pp. 2078–2089, 2018.
- [15] Q. Jiang, F. Shao, W. Gao, Z. Chen, G. Jiang, and Y.-S. Ho, "Unified no-reference quality assessment of singly and multiply distorted stereoscopic images," *IEEE Transactions on Image Processing*, vol. 28, no. 4, pp. 1866–1881, 2019.
- [16] Q. Jiang, W. Gao, S. Wang, G. Yue, F. Shao, Y.-S. Ho, and S. Kwong, "Blind image quality measurement by exploiting high-order statistics with deep dictionary encoding network," *IEEE Transactions on Instrumentation and Measurement*, vol. 69, no. 10, pp. 7398–7410, 2020.
- [17] L. Li, Y. Zhou, J. Wu, F. Li, and G. Shi, "Quality index for view synthesis by measuring instance degradation and global appearance," *IEEE Transactions on Multimedia*, vol. 23, pp. 320–332, 2021.
- [18] M. Yu, H. Lakshman, and B. Girod, "A framework to evaluate omnidirectional video coding schemes," in *2015 IEEE International Symposium on Mixed and Augmented Reality*, 2015, pp. 31–36.
- [19] Y. Zhou, Y. Ding, Y. Sun, L. Li, J. Wu, and X. Gao, "Perceptual information completion-based siamese omnidirectional image quality assessment network," *IEEE Transactions on Instrumentation and Measurement*, vol. 73, pp. 1–10, 2024.
- [20] J. Xu, W. Zhou, and Z. Chen, "Blind omnidirectional image quality assessment with viewport oriented graph convolutional networks," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 5, pp. 1724–1737, 2021.
- [21] Y. Liu, X. Yin, Y. Wang, Z. Yin, and Z. Zheng, "Hvs-based perception-driven no-reference omnidirectional image quality assessment," *IEEE Transactions on Instrumentation and Measurement*, vol. 72, pp. 1–11, 2023.
- [22] V. Zakharchenko, K. P. Choi, and J. H. Park, "Quality metric for spherical panoramic video," in *Optics and Photonics for Information Processing X*, K. M. Iftekharuddin, A. A. S. Awwal, M. G. Vázquez, A. Márquez, and M. A. Matin, Eds., vol. 9970, International Society for Optics and Photonics. SPIE, 2016, p. 99700C. [Online]. Available: <https://doi.org/10.1117/12.2235885>
- [23] Y. Zhang and D. M. Chandler, "An algorithm for no-reference image quality assessment based on log-derivative statistics of natural scenes," in *Image Quality and System Performance X*, P. D. Burns and S. Triantaphillidou, Eds., vol. 8653, International Society for Optics and Photonics. SPIE, 2013, p. 86530J. [Online]. Available: <https://doi.org/10.1117/12.2001342>
- [24] P. C. Madhusudana and R. Soundararajan, "Subjective and objective quality assessment of stitched images for virtual reality," *IEEE Transactions on Image Processing*, vol. 28, no. 11, pp. 5620–5635, 2019.

- [25] K. Ma, W. Liu, T. Liu, Z. Wang, and D. Tao, "dipiQ: Blind image quality assessment by learning-to-rank discriminable image pairs," *IEEE Transactions on Image Processing*, vol. 26, no. 8, pp. 3951–3964, 2017.
- [26] 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, 2018.
- [27] H. Jiang, G. Jiang, M. Yu, Y. Zhang, Y. Yang, Z. Peng, F. Chen, and Q. Zhang, "Cubemap-based perception-driven blind quality assessment for 360-degree images," *IEEE Transactions on Image Processing*, vol. 30, pp. 2364–2377, 2021.
- [28] Y. Liu, J. Wu, L. Li, W. Dong, and G. Shi, "Quality assessment of ugc videos based on decomposition and recomposition," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 33, no. 3, pp. 1043–1054, 2023.
- [29] W. Sun, W. Luo, X. Min, G. Zhai, X. Yang, K. Gu, and S. Ma, "Mc360iqa: The multi-channel cnn for blind 360-degree image quality assessment," in *2019 IEEE International Symposium on Circuits and Systems (ISCAS)*, 2019, pp. 1–5.
- [30] X. Min, G. Zhai, K. Gu, Y. Liu, and X. Yang, "Blind image quality estimation via distortion aggravation," *IEEE Transactions on Broadcasting*, vol. 64, no. 2, pp. 508–517, 2018.
- [31] R. Xu, S. Yang, Y. Wang, B. Du, and H. Chen, "A survey on vision mamba: Models, applications and challenges," 2024.
- [32] Y. Sun, A. Lu, and L. Yu, "Weighted-to-spherically-uniform quality evaluation for omnidirectional video," *IEEE Signal Processing Letters*, vol. 24, no. 9, pp. 1408–1412, 2017.
- [33] S. Chen, Y. Zhang, Y. Li, Z. Chen, and Z. Wang, "Spherical structural similarity index for objective omnidirectional video quality assessment," in *2018 IEEE International Conference on Multimedia and Expo (ICME)*, 2018, pp. 1–6.
- [34] H. G. Kim, H.-T. Lim, and Y. M. Ro, "Deep virtual reality image quality assessment with human perception guider for omnidirectional image," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 4, pp. 917–928, 2020.
- [35] Y. Liu, H. Yu, B. Huang, G. Yue, and B. Song, "Blind omnidirectional image quality assessment based on structure and natural features," *IEEE Transactions on Instrumentation and Measurement*, vol. 70, pp. 1–11, 2021.
- [36] Y. Zhou, Y. Sun, L. Li, K. Gu, and Y. Fang, "Omnidirectional image quality assessment by distortion discrimination assisted multi-stream network," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 4, pp. 1767–1777, 2022.
- [37] A. Gu and T. Dao, "Mamba: Linear-time sequence modeling with selective state spaces," 2024.
- [38] A. Gu, K. Goel, and C. Ré, "Efficiently modeling long sequences with structured state spaces," 2022.
- [39] A. Gu, A. Gupta, K. Goel, and C. Ré, "On the parameterization and initialization of diagonal state space models," in *Proceedings of the 36th International Conference on Neural Information Processing Systems*, ser. NIPS '22. Red Hook, NY, USA: Curran Associates Inc., 2024.
- [40] A. Gupta, A. Gu, and J. Berant, "Diagonal state spaces are as effective as structured state spaces," in *Proceedings of the 36th International Conference on Neural Information Processing Systems*, ser. NIPS '22. Red Hook, NY, USA: Curran Associates Inc., 2024.
- [41] A. Orvieto, S. L. Smith, A. Gu, A. Fernando, C. Gulcehre, R. Pascanu, and S. De, "Resurrecting recurrent neural networks for long sequences," in *Proceedings of the 40th International Conference on Machine Learning*, ser. ICML'23. JMLR.org, 2023.
- [42] J. T. H. Smith, A. Warrington, and S. W. Linderman, "Simplified state space layers for sequence modeling," 2023.
- [43] R. Hasani, M. Lechner, T.-H. Wang, M. Chahine, A. Amini, and D. Rus, "Liquid structural state-space models," 2022.
- [44] A. Gu, I. Johnson, A. Timalina, A. Rudra, and C. Ré, "How to train your hippo: State space models with generalized orthogonal basis projections," 2022.
- [45] J. Wang, W. Zhu, P. Wang, X. Yu, L. Liu, M. Omar, and R. Hamid, "Selective structured state-spaces for long-form video understanding," in *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 6387–6397.
- [46] H. Mehta, A. Gupta, A. Cutkosky, and B. Neyshabur, "Long range language modeling via gated state spaces," 2022.
- [47] A. Chubarau and J. Clark, "Vtamiq: Transformers for attention modulated image quality assessment," 2021.
- [48] Y. Liu, Y. Tian, Y. Zhao, H. Yu, L. Xie, Y. Wang, Q. Ye, and Y. Liu, "Vmamba: Visual state space model," 2024.
- [49] Z. Liu, H. Hu, Y. Lin, Z. Yao, Z. Xie, Y. Wei, J. Ning, Y. Cao, Z. Zhang, L. Dong, F. Wei, and B. Guo, "Swin transformer v2: Scaling up capacity and resolution," in *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 11 999–12 009.
- [50] C.-S. Chen, G.-Y. Chen, D. Zhou, D. Jiang, and D.-S. Chen, "Resvmamba: Fine-grained food category visual classification using selective state space models with deep residual learning," 2024.
- [51] C. Du, Y. Li, and C. Xu, "Understanding robustness of visual state space models for image classification," 2024.
- [52] Y. Shi, B. Xia, X. Jin, X. Wang, T. Zhao, X. Xia, X. Xiao, and W. Yang, "Vmambair: Visual state space model for image restoration," 2024.
- [53] Z. Wang, J.-Q. Zheng, C. Ma, and T. Guo, "Vmambamorph: a multi-modality deformable image registration framework based on visual state space model with cross-scan module," 2024.
- [54] H.-Y. Ma, L. Zhang, and S. Shi, "Vmambacc: A visual state space model for crowd counting," 2024.
- [55] X. Xie, Y. Cui, C.-I. Jeong, T. Tan, X. Zhang, X. Zheng, and Z. Yu, "Fusionmamba: Dynamic feature enhancement for multimodal image fusion with mamba," 2024.
- [56] Y. Yang, C. Ma, J. Yao, Z. Zhong, Y. Zhang, and Y. Wang, "Remember: Referring image segmentation with mamba twister," 2024.
- [57] C. Tian, F. Shao, X. Chai, Q. Jiang, L. Xu, and Y.-S. Ho, "Viewportsphere-branch network for blind quality assessment of stitched 360° omnidirectional images," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 33, no. 6, pp. 2546–2560, 2023.
- [58] T. T. Huong, D. T. Ha, H. T. T. Tran, N. D. Viet, B. D. Tien, N. H. Thanh, T. C. Thang, and P. N. Nam, "An effective foveated 360° image assessment based on graph convolution network," *IEEE Access*, vol. 10, pp. 98 165–98 178, 2022.
- [59] J. Hu, X. Wang, F. Shao, and Q. Jiang, "Tspr: Deep network-based blind image quality assessment using two-side pseudo reference images," *Digital Signal Processing*, vol. 106, p. 102849, 2020. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1051200420301949>
- [60] M. V. Conde, G. Geigle, and R. Timofte, "Instructir: High-quality image restoration following human instructions," 2024.
- [61] T. Brooks, A. Holynski, and A. A. Efros, "Instructpix2pix: Learning to follow image editing instructions," in *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 18 392–18 402.
- [62] K. Zhang, J. Liang, L. Van Gool, and R. Timofte, "Designing a practical degradation model for deep blind image super-resolution," in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 4771–4780.
- [63] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [64] X. Pei, T. Huang, and C. Xu, "Efficientvmamba: Atrous selective scan for light weight visual mamba," 2024.
- [65] Y. Fan and C. Chen, "Omiquet: Multiscale feature aggregation convolutional neural network for omnidirectional image assessment," *Applied Intelligence*, vol. 54, no. 7, pp. 5711–5727, Apr 2024. [Online]. Available: <https://doi.org/10.1007/s10489-024-05421-1>
- [66] H. Duan, G. Zhai, X. Min, Y. Zhu, Y. Fang, and X. Yang, "Perceptual quality assessment of omnidirectional images," in *2018 IEEE International Symposium on Circuits and Systems (ISCAS)*, 2018, pp. 1–5.
- [67] G. M. Pace, M. T. Ivancic, G. L. Edwards, B. A. Iwata, and T. J. Page, "Assessment of stimulus preference and reinforcer value with profoundly retarded individuals," *J Appl Behav Anal*, vol. 18, no. 3, pp. 249–255, 1985.
- [68] G. Wallace, "The jpeg still picture compression standard," *IEEE Transactions on Consumer Electronics*, vol. 38, no. 1, pp. xviii–xxxiv, 1992.
- [69] T. Wiegand, G. Sullivan, G. Bjontegaard, and A. Luthra, "Overview of the h.264/avc video coding standard," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 13, no. 7, pp. 560–576, 2003.
- [70] G. J. Sullivan, J.-R. Ohm, W.-J. Han, and T. Wiegand, "Overview of the high efficiency video coding (hevc) standard," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 22, no. 12, pp. 1649–1668, 2012.
- [71] H. Duan, X. Min, W. Sun, Y. Zhu, X.-P. Zhang, and G. Zhai, "Attentive deep image quality assessment for omnidirectional stitching," *IEEE Journal of Selected Topics in Signal Processing*, vol. 17, no. 6, pp. 1150–1164, 2023.
- [72] L. Li, Y. Zhou, W. Lin, J. Wu, X. Zhang, and B. Chen, "No-reference quality assessment of deblocked images," *Neurocomput.*, vol. 177, no. C, p. 572–584, feb 2016. [Online]. Available: <https://doi.org/10.1016/j.neucom.2015.11.063>

- [73] X. Min, K. Ma, K. Gu, G. Zhai, Z. Wang, and W. Lin, "Unified blind quality assessment of compressed natural, graphic, and screen content images," *IEEE Transactions on Image Processing*, vol. 26, no. 11, pp. 5462–5474, 2017.
- [74] A. Mittal, A. K. Moorthy, and A. C. Bovik, "No-reference image quality assessment in the spatial domain," *IEEE Transactions on Image Processing*, vol. 21, no. 12, pp. 4695–4708, 2012.
- [75] L. Li, Y. Zhou, K. Gu, W. Lin, and S. Wang, "Quality assessment of dibr-synthesized images by measuring local geometric distortions and global sharpness," *IEEE Transactions on Multimedia*, vol. 20, no. 4, pp. 914–926, 2018.
- [76] L. Li, Y. Zhou, K. Gu, Y. Yang, and Y. Fang, "Blind realistic blur assessment based on discrepancy learning," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 11, pp. 3859–3869, 2020.
- [77] C. Li, M. Xu, X. Du, and Z. Wang, "Bridge the gap between vqa and human behavior on omnidirectional video: A large-scale dataset and a deep learning model," in *Proceedings of the 26th ACM International Conference on Multimedia*, ser. MM '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 932–940. [Online]. Available: <https://doi.org/10.1145/3240508.3240581>
- [78] C. Li, M. Xu, L. Jiang, S. Zhang, and X. Tao, "Viewport proposal cnn for 360° video quality assessment," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 10169–10178.
- [79] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, *PyTorch: an imperative style, high-performance deep learning library*. Red Hook, NY, USA: Curran Associates Inc., 2019.
- [80] H. Robbins and S. Monro, "A Stochastic Approximation Method," *The Annals of Mathematical Statistics*, vol. 22, no. 3, pp. 400 – 407, 1951. [Online]. Available: <https://doi.org/10.1214/aoms/1177729586>
- [81] J. Kiefer and J. Wolfowitz, "Stochastic Estimation of the Maximum of a Regression Function," *The Annals of Mathematical Statistics*, vol. 23, no. 3, pp. 462 – 466, 1952. [Online]. Available: <https://doi.org/10.1214/aoms/1177729392>
- [82] S. Ruder, "An overview of gradient descent optimization algorithms," 2017.
- [83] 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, 2013.
- [84] P. Ye, J. Kumar, L. Kang, and D. Doermann, "Unsupervised feature learning framework for no-reference image quality assessment," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 1098–1105.
- [85] W. Xue, L. Zhang, and X. Mou, "Learning without human scores for blind image quality assessment," in *2013 IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 995–1002.
- [86] 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.
- [87] Q. Wu, Z. Wang, and H. Li, "A highly efficient method for blind image quality assessment," in *2015 IEEE International Conference on Image Processing (ICIP)*, 2015, pp. 339–343.
- [88] J. Xu, P. Ye, Q. Li, H. Du, Y. Liu, and D. Doermann, "Blind image quality assessment based on high order statistics aggregation," *IEEE Transactions on Image Processing*, vol. 25, no. 9, pp. 4444–4457, 2016.
- [89] X. Min, K. Gu, G. Zhai, J. Liu, X. Yang, and C. W. Chen, "Blind quality assessment based on pseudo-reference image," *IEEE Transactions on Multimedia*, vol. 20, no. 8, pp. 2049–2062, 2018.



**Yun Liu** received the Ph.D. degree in communication and information engineering from Tianjin University, China, in 2016. From 2014 to 2015, she was a visiting Ph.D. student at the Visual Space Perception Laboratory, University of California, Berkeley, United States.

She is currently an associate professor at the Faculty of Information, Liaoning University, Shenyang, China. Her research interests include multimedia quality assessment, image processing, computer vision, and pattern recognition.



**Sifan Li** is currently pursuing the M.S. degree in computer science and technology at the Faculty of Information, Liaoning University, Shenyang, China. His research interests include multimedia quality assessment, image processing, efficient training and inference, and computer vision.



**Huiyu Duan** received the B.E. degree from the University of Electronic Science and Technology of China, Chengdu, China, in 2017, and the Ph.D. degree from Shanghai Jiao Tong University, Shanghai, China, in 2024. He is currently a Postdoctoral Fellow at Shanghai Jiao Tong University. From Sept. 2019 to Sept. 2020, he was a visiting Ph.D. student at the Schepens Eye Research Institute, Harvard Medical School, Boston, USA. He received the Best Paper Award of IEEE International Symposium on Broadband Multimedia Systems and Broadcasting (BMSB) in 2022. His research interests include perceptual quality assessment, quality of experience, visual attention modeling, extended reality (XR), and multimedia signal processing.



**Yu Zhou** received the B.S. and Ph.D. degrees from the China University of Mining and Technology, Xuzhou, China, in 2014 and 2019, respectively.

She is currently an Associate Professor with the School of Information and Control Engineering, China University of Mining and Technology. Her research interests include computer vision, multimedia image processing, and artificial intelligence.



**Daoxin Fan** is currently pursuing the M.S. degree in computer science and technology at the Faculty of Information, Liaoning University, Shenyang, China. His research interests include multimedia quality assessment, image processing, and computer vision.



**Guangtao Zhai** (Senior Member, IEEE) received the B.E. and M.E. degrees from Shandong University, Shandong, China, in 2001 and 2004, respectively, and the Ph.D. degree from Shanghai Jiao Tong University, Shanghai, China, in 2009, where he is currently a Research Professor with the Institute of Image Communication and Information Processing. From 2008 to 2009, he was a Visiting Student with the Department of Electrical and Computer Engineering, McMaster University, Hamilton, ON, Canada, where he was a Post-Doctoral Fellow from 2010 to 2012. From 2012 to 2013, he was a Humboldt Research Fellow with the Institute of Multimedia Communication and Signal Processing, Friedrich Alexander University of Erlangen-Nuremberg, Germany. He received the Award of National Excellent Ph.D. Thesis from the Ministry of Education of China in 2012. His research interests include multimedia signal processing and perceptual signal processing.