

# AIM 2024 Challenge on UHD Blind Photo Quality Assessment

Vlad Hosu<sup>3†</sup>, Marcos V. Conde<sup>1,2†</sup>, Lorenzo Agnolucci<sup>3,4†</sup>, Nabajeet Barman<sup>2†</sup>, Saman Zadtootaghaj<sup>2†</sup>, Radu Timofte<sup>1†</sup>, Wei Sun, Weixia Zhang, Yuqin Cao, Linhan Cao, Jun Jia, Zijian Chen, Zicheng Zhang, Xionghuo Min, Guangtao Zhai, Songbai Tan, Lixin Zhang, Guanghui Yue, Daekyu Kwon, Dongyoung Kim, Seon Joo Kim, Yunchen Zhang, Xiangkai Xu, Hong Gao, Yiming Bao, Ji Shi, Xiugang Dong, Xiangsheng Zhou, Yaofeng Tu, Zewen Chen, Shunhan Xu, Haochen Guo, Yun Zeng, Shuai Liu, Jian Guo, Juan Wang, Bing Li, Dehua Liu, Hesong Liu, Grigory Malivenko, Asile Gerek, Xingyuan Ma, Cheng Li, Joonhee Lee, Junseo Bang, and Se Young Chun

<sup>1</sup> Computer Vision Lab, CAIDAS & IFI, University of Würzburg

<sup>2</sup> Visual Computing Group, FTG, Sony PlayStation

<sup>3</sup> Sony AI

<sup>4</sup> University of Florence

† Challenge Organizers, ‡ Corresponding Author

<https://database.mmsp-kn.de/uhd-iqa-benchmark-database.html>



**Fig. 1:** Example images from the UHD-IQA dataset [14]. They have been cropped to 64% of their original size to enhance detail visibility. The author's name from [Pixabay.com](https://www.pixabay.com) is shown at the bottom right of each image.

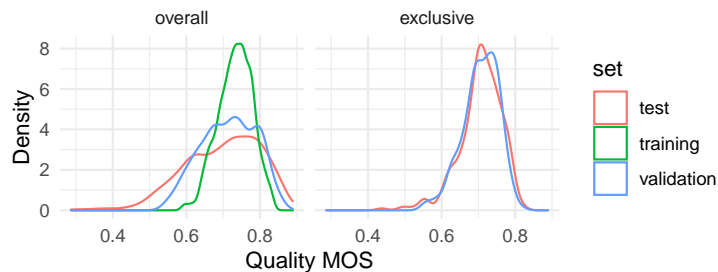
**Abstract.** We introduce the AIM 2024 UHD-IQA Challenge, a competition to advance the No-Reference Image Quality Assessment (NR-IQA) task for modern, high-resolution photos. The challenge is based on the recently released UHD-IQA Benchmark Database, which comprises 6,073 UHD-1 (4K) images annotated with perceptual quality ratings from expert raters. Unlike previous NR-IQA datasets, UHD-IQA focuses on highly aesthetic photos of superior technical quality, reflecting the ever-increasing standards of digital photography. This challenge aims to develop efficient and effective NR-IQA models. Participants are tasked with creating novel architectures and training strategies to achieve high predictive performance on UHD-1 images within a computational budget of 50G MACs. This enables model deployment on edge devices and scalable processing of extensive image collections. Winners are determined based on a combination of performance metrics, including correlation measures (SRCC, PLCC, KRCC), absolute error metrics (MAE, RMSE), and computational efficiency (G MACs). To excel in this challenge, participants leverage techniques like knowledge distillation, low-precision inference, and multi-scale training. By pushing the boundaries of NR-IQA for high-resolution photos, the UHD-IQA Challenge aims to stimulate the development of practical models that can keep pace with the rapidly evolving landscape of digital photography. The innovative solutions emerging from this competition will have implications for various applications, from photo curation and enhancement to image compression.

## 1 Introduction

Blind Image Quality Assessment (BIQA) is essential for various applications, including camera benchmarking, professional photo curation, and image enhancement. Despite advances in BIQA models, their effectiveness is constrained by the limitations of existing datasets. Current datasets are primarily annotated at standard definition (SD) resolutions and focus on images with obvious distortions. As a result, BIQA models struggle with high-resolution images that exhibit subtle degradations, which are increasingly common with modern cameras.

These datasets also suffer from a bias toward average or low-quality images, leading to a class imbalance that weakens the generalization of BIQA models. As camera technology advances, producing higher-quality and higher-resolution images, the need for better datasets becomes critical. Moreover, the efficient processing of these high-quality images on edge devices or at scale remains challenging, as most current models are not optimized for such tasks.

We introduce the UHD-IQA challenge as part of AIM 2024 to address these issues. The UHD-IQA benchmark dataset focuses on ultra-high-definition (UHD) images of high aesthetic and technical quality, aiming to fill the gaps in existing benchmarks. The challenge is developing efficient BIQA models that fully leverage this dataset, ensuring high accuracy and computational efficiency for real-world applications.



**Fig. 2:** Density of quality MOS per subset. "Overall" includes all image categories, whereas "exclusive" refers to categories that are only part of the validation and test sets.

### 1.1 UHD-IQA Benchmark Database

The dataset comprises 6073 ultra-high-definition (UHD-1, 4K) images, all annotated at a fixed width of 3840 pixels. Unlike existing BIQA datasets, ours focuses on high-quality images with a strong aesthetic appeal, filling a critical literature gap. The images were sourced from Pixabay.com, a repository of CC0-licensed stock photos, and were manually curated to exclude synthetic or heavily edited content. This ensures that the dataset consists of genuine, high-quality photographs. The dataset split is as follows: 4269 for training, 904 for validation, and 900 for testing.

We conducted a crowdsourcing study involving ten expert raters, including photographers and graphic artists, to achieve reliable annotations. Each expert assessed each image at least twice in multiple sessions, yielding 20 ratings per image. The rigorous annotation process and rich metadata, including user and machine-generated tags from over 5,000 categories, provide a comprehensive and reliable resource for training BIQA models.

Furthermore, the test and validation sets include a special subset of 300 images out of approximately 900 in each set, labeled as "*exclusive*" – see the MOS density in Fig. 2. This subset is selected based on image categories excluded from the training set. The categories for all images were either automatically annotated using AWS Rekognition or manually specified by the image authors when publishing to [Pixabay.com](https://pixabay.com).

The exclusive categories were chosen to be distinct from typical ImageNet ones, focusing on images that do not feature a single dominant object. Instead, they depict multiple scattered objects or wide-spanning scenes. This selection aims to encourage the use of more general-purpose pre-training features. The exclusive categories are Sea, Ocean, Sand, Landscape, Mountain(s), Scenery, City, and Urban.

The performance on the exclusive split also provides valuable insights into each model's generalization capabilities when deviating from the image distribution of the training set.

## 1.2 The AIM 2024 Challenge

The challenge participants were tasked with developing novel BIQA models that efficiently and effectively assess high-resolution images. The proposed models were required to operate below 50 GMACs, ensuring they are lightweight enough for deployment on edge devices or scalable processing. Participants were encouraged to employ strategies such as knowledge distillation and low-precision inference and to select optimal pre-training datasets to meet these requirements.

The challenge was structured around multiple evaluation criteria to determine individual rankings. These criteria included correlation metrics – Pearson Linear Correlation Coefficient (PLCC), Spearman Rank-order Correlation Coefficient (SRCC), and Kendall Rank Correlation Coefficient (KRCC) – as well as absolute error metrics such as Mean Absolute Error (MAE) and Root Mean Square Error (RMSE). Additionally, compute efficiency was a critical factor in determining the winning models. By pushing the boundaries of BIQA with this challenge, we aim to drive the development of practical, scalable, and high-performing models that are well-suited for modern, high-quality images.

*Associated AIM Challenges.* This challenge is one of the AIM 2024 Workshop<sup>5</sup> associated challenges on: sparse neural rendering [28,29], UHD blind photo quality assessment [15], compressed depth map super-resolution and restoration [11], efficient video super-resolution for AV1 compressed content [10], video super-resolution quality assessment [25], compressed video quality assessment [33] and video saliency prediction [26].

## 2 Proposed Methods

Eight methods were submitted for the final round of the challenge. Most solutions consist of ensembles of multiple neural networks, especially Transformer-based [13,24] models and CLIP-based [20] models.

As a *baseline*, we propose an efficient solution based on MobileNet V3 [17,18]. The original high-resolution images are cropped (focusing on the center) at  $960 \times 1920$  pixels; these are resized to HD resolution ( $1280 \times 720$ ). Using a fine-tuned MobileNet V3 [17] backbone as a feature extractor allows to reduce overfitting and training time and faster inference speed. The baseline model has 3.22 M parameters and a computational complexity of 4.2 GMACs.

### 2.1 Challenge Results

Table 1 and Table 2 present comparative evaluation results of the eight teams' performance in predicting the quality MOS using various metrics.

The top three performances for each metric are highlighted, with gold, silver, and bronze representing the first, second, and third-best results, respectively.

<sup>5</sup> <https://www.cvlai.net/aim/2024/>

Method	MAE ↓	RMSE ↓	PLCC ↑	SRCC ↑	KRCC ↑
SJTU (2.2)	0.0418	0.0615	0.7985	0.8463	0.6573
GS-PIQA (2.3)	0.0430	0.0607	0.7925	0.8297	0.6399
CIPLAB (2.4)	0.0445	0.0638	0.7995	0.8354	0.6419
EQCNet (2.5)	0.0438	0.0621	0.7682	0.7954	0.6055
MobileNet-IQA (2.6)	0.0463	0.0659	0.7558	0.7883	0.5975
NF-RegNets (2.7)	0.0494	0.0703	0.7222	0.7715	0.5806
Challenge Baseline	0.0502	0.0733	0.6881	0.7462	0.5537
CLIP-IQA* (2.8)	0.0519	0.0723	0.7116	0.7305	0.5393
ICL (2.9)	0.1147	0.1364	0.5206	0.5166	0.3615
HyperIQA [34]	0.070	0.118	0.103	0.553	0.389
Effnet-2C-MLSP [42]	0.059	0.074	0.641	0.675	0.491
CONTRIQUE [23]	0.052	0.073	0.678	0.732	0.532
ARNIQA [2]	0.052	0.074	0.694	0.739	0.544
CLIP-IQA+ [40]	0.089	0.111	0.709	0.747	0.551
QualiCLIP [1]	0.066	0.083	0.725	0.770	0.570

**Table 1:** Official test split performance. We highlight the top-3 (gold, silver, bronze) methods for the different metrics. The top section lists methods that participated in the AIM 2024 challenge. The bottom section presents baselines derived from retraining existing methods, which require more than 200 GMACs.

However, the winner and runner-up teams are ranked considering the final score for each team, which is computed as follows.

Let  $\mathcal{S}_i$  denote the main score for team  $i$ , and  $\mathcal{R}(\mathcal{M}_i) = \text{Rank}(\mathcal{M}_i)$ ,  $\mathcal{R}(\mathcal{M}_i) = 1, \dots, N$  is the ranking function that assigns a rank based on the metric value  $\mathcal{M}_i$  for each of the  $N = 8$  participating teams. The best rank is 1. Correlation metrics are ranked highest when they have higher values, whereas absolute errors rank best when they are lowest.

$$\mathcal{S}_i = \frac{1}{5} \left[ \mathcal{R}(\mathcal{M}_i^{\text{MAE}}) + \mathcal{R}(\mathcal{M}_i^{\text{RMSE}}) + \mathcal{R}(\mathcal{M}_i^{\text{KRCC}}) + \mathcal{R}(\mathcal{M}_i^{\text{PLCC}}) + \mathcal{R}(\mathcal{M}_i^{\text{SRCC}}) \right]$$

$\mathcal{M}_i^{\text{Metric}}$  represents the value of the previously mentioned metrics.

The team with the lowest main score  $\mathcal{S}_i$  is considered to be the winner. Based on the scores obtained and shown in Table 2, team SJTU is the overall competition winner, followed by team SZU SongBai (first runner-up) and team CIPLAB (second runner-up).

Table 3 presents a comparative evaluation of the teams' performance in predicting MOS across various evaluation metrics specifically on the exclusive portion of the test set. As expected, the results indicate a noticeable reduction in performance metrics. Interestingly, CIPLAB ranks second in this evaluation (compared to third in the overall ranking in Table 1), which might be due to better generalization capabilities of the model compared to GS-PIQA.

Models	MAE ↓	RMSE ↓	PLCC ↑	SRCC ↑	KRCC ↑
EQCNet (2.5)	0.0299	0.0383	0.8285	0.8234	0.6342
SJTU (2.2)	0.0318	0.0402	0.8238	0.8169	0.6244
GS-PIQA (2.3)	0.0332	0.0406	0.8192	0.8092	0.6181
CIPLAB (2.4)	0.0329	0.0423	0.8136	0.8063	0.6143
MobileNet-IQA (2.6)	0.0345	0.0439	0.7831	0.7757	0.5824
NF-RegNets (2.7)	0.0352	0.0444	0.7968	0.7897	0.5973
Challenge Baseline	0.0372	0.0482	0.7445	0.7422	0.5504
CLIP-IQA* (2.8)	0.0398	0.0509	0.7069	0.6918	0.5112
ICL (2.9)	0.0622	0.0737	0.5217	0.5101	0.3580
HyperIQA [34]	0.055	0.087	0.182	0.524	0.359
Effnet-2C-MLSP [42]	0.050	0.060	0.627	0.615	0.445
CONTRIQUE [23]	0.038	0.049	0.712	0.716	0.521
ARNIQA [2]	0.039	0.050	0.717	0.718	0.523
CLIP-IQA+ [40]	0.087	0.108	0.732	0.743	0.546
QualiCLIP [1]	0.064	0.079	0.752	0.757	0.557

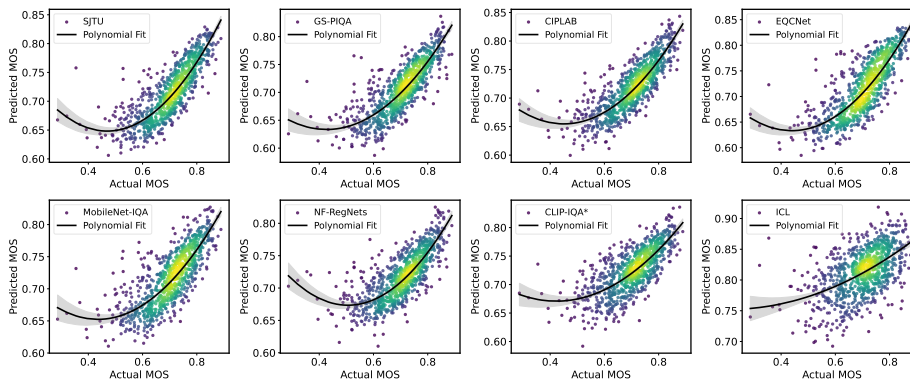
**Table 2:** Official **validation** split performance. Comparison of models with top-3 (gold, silver, bronze) highlighted for each metric. The top section lists methods that participated in the AIM 2024 challenge. The bottom section presents baselines derived from retraining existing methods, which require more than 200 GMACs.

Method	MAE ↓	RMSE ↓	PLCC ↑	SRCC ↑	KRCC ↑
SJTU (2.2)	0.0292	0.0422	0.6816	0.7407	0.5471
CIPLAB (2.4)	0.0308	0.0439	0.6733	0.7009	0.5078
GS-PIQA (2.3)	0.0320	0.0447	0.6325	0.6710	0.4915
EQCNet (2.5)	0.0328	0.0453	0.6227	0.6555	0.4786
MobileNet-IQA (2.6)	0.0328	0.0466	0.5916	0.5999	0.4320
NF-RegNets (2.7)	0.0338	0.0480	0.5707	0.6099	0.4388
CLIP-IQA* (2.8)	0.0361	0.0510	0.5113	0.5157	0.3622
ICL (2.9)	0.1014	0.1138	0.4331	0.4106	0.2802

**Table 3:** The performance evaluation of **exclusive** test split. We highlight the top-3 (gold, silver, bronze) methods for the different metrics.

Figure 3 presents a comparative analysis of predicted quality scores against ground-truth MOS for the eight competing teams. Each subplot represents the performance of a particular team, with the team name shown in the legend. The X-axis values represent ground-truth (actual) MOS; the predicted scores are shown on the y-axis. The purple scatter points represent a particular image prediction score, with higher-density areas shown in yellow. The polynomial fit, shown as a black curve, highlights the general trend in the predictions relative to the ground truth.

It can be observed that all teams display a positive correlation between the predicted and ground-truth MOS, as indicated by the upward trend in all subplots. However, the digress of scatter around the fitted curve varies across the subplots indicating the difference in the strength and alignment of this correlation between various teams. For example, the polynomial fit for teams like



**Fig. 3:** Scatter plots of the predicted quality scores vs ground-truth (actual) MOS. The curves were obtained by a second-order polynomial fitting.

‘*SJTU (2.2)*’, ‘*GS-PIQA (2.3)*’ and ‘*CIPLAB (2.4)*’ show a tighter clustering of data points around the curve, which indicates a better alignment of predicted quality with the ground-truth MOS. On the other hand, teams such as ‘*ICL (ef-sec:icl)*’ and ‘*NF-RegNets (2.7)*’ show more scattered data points. Overall, while all teams demonstrate the ability to predict MOS scores with some degree of accuracy, there are clear differences in prediction quality.

Method	Input	Training Time (hrs)	Extra Data	Params. (M)	MACs (G)	GPU
SJTU (2.2)	$480 \times 480$	12	Yes	82.85	43.53	RTX 3090
GS-PIQA (2.3)	$384 \times 384$	4	No	144.814	50.260	GTX3090
CIPLAB (2.4)	$2160 \times 3840$	12	No	113	44	RTX 2080 Ti
EQCNet (2.5)	$384 \times 384 - 1366 \times 768$	22	Yes	30.15	12.97	A800
MobileViT-IQA (2.6)	$1907 \times 1231$	18	No	96.72	359.74	A800
MobileNet-IQA (2.6)	$1907 \times 1231$	48	No	81.48	46.73	A800
NF-RegNets (2.7)	$720 \times 720$	$\approx 10$	No	28.5	44.52	$2 \times 2070$ Ti
CLIP-IQA* (2.8)	$224 \times 224$	0.25	No	151	48.5	A6000
ICL (2.9)	$2160 \times 3840$	0.1	No	139.1	42.09	A100
Challenge Baseline	$1280 \times 720$	6	No	3.2	4.2	3090Ti

**Table 4:** Training specification for each method. All inputs are 3-channel RGB images; only the spatial dimensions are listed.

**Summary of Implementation Details** A summary of the methods is provided in Table 4, which includes details on the input resolution, computational complexity measured in MACs, and the number of parameters for each model.

In the following sections, we describe the top solutions to the challenge. Please note that the method descriptions were provided by the respective teams or individual participants as their contributions to this report.

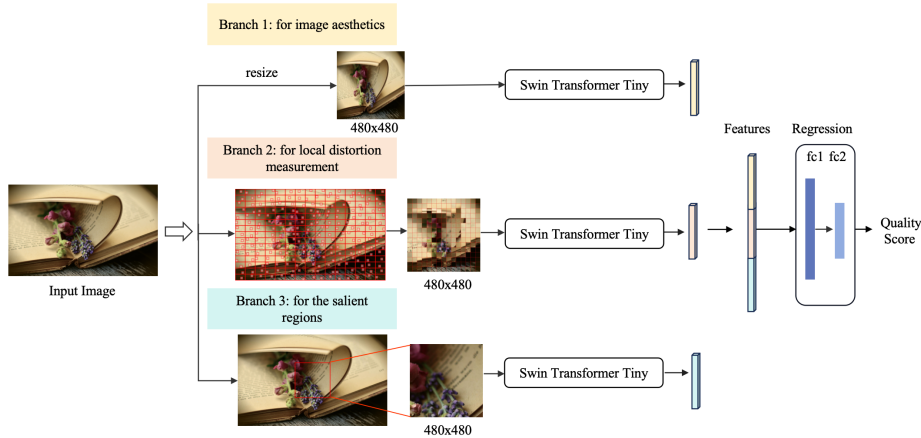


Fig. 4: The method proposed by the SJTU Team, using three branches [37].

## 2.2 Assessing UHD Image Quality from Aesthetics, Distortion, and Saliency

*Wei Sun, Weixia Zhang, Yuqin Cao, Linhan Cao, Jun Jia, Zijian Chen, Zicheng Zhang, Xionghuo Min, Guangtao Zhai  
Shanghai Jiao Tong University (SJTU), China*

We design a multi-branch deep neural network (DNN) to evaluate the UHD image quality from three perspectives: **global aesthetic characteristics, local technique distortions, and salient region perception**, while avoiding direct processing of high-resolution images [37]. Specifically, a low-resolution image resized from the UHD image, a fragment image composed of local fragments cropped from the equal-size patches of the UHD image, and the center patch cropped from the UHD image are used as inputs to extract the respective features through three branches. The Swin Transformer Tiny [22] pre-trained on the AVA dataset [27] are utilized as the backbone networks of the three branches. The extracted features are concatenated and regressed into quality scores by a two-layer multi-layer perceptron (MLP). We employ the mean square error (MSE) loss and the fidelity loss [39] to optimize the proposed model. By dividing the overall quality measurement of the high-resolution image into three quality dimension measurements of low-resolution images, our method effectively assesses the quality of UHD images with an acceptable computational complexity. Moreover, we avoid complex model designs and use only the standard DNN structures, making it easy to implement in practical applications and optimize for hardware.

The proposed model is illustrated in Fig. 4. It consists of three branches to extract the quality-aware features from aspects of global aesthetic characteristics, local technique distortions, and salient object perception.



First, we consider **image aesthetics**, which encompasses the overall perception of image characteristics such as content, layout, color, contrast, etc. These usually are global features that do not require high resolution. Thus, we resize the UHD image to a low resolution of  $480 \times 480$  and use the low-resolution image as the input of the branch responsible for the aesthetic characteristics.

Second, we address **low-level image distortions**, which are typically evident on local image patches and are sensitive to the resolution. We employ a fragment sampling strategy [43], where the entire image is divided into  $15 \times 15$  equal-sized patches, and a smaller fragment with a resolution of  $32 \times 32$  is randomly cropped from each patch. These fragments are then spliced into a fragment image of  $480 \times 480$ , which serves as input to the branch responsible for local distortion measurement.

Third, since UHD images are often viewed on large screens where the human visual system tends to focus on salient regions, **the quality of the salient region is crucial for the overall quality**. Considering the center bias of saliency detection [3], we crop the center patch with a resolution of  $480 \times 480$  from the UHD image to extract the quality-aware features for the salient regions.

Finally, we use Swin Transformer Tiny [22] pre-trained on the AVA dataset as the backbone of three branches to extract the corresponding features for each aspect. Note that these three branches do not share the model weights. The extracted features are concatenated as the quality-aware feature representation and then regressed into quality scores via a two-layer MLP network. The two-layer MLP network consists of 128 and 1 neurons, respectively. We employ the mean square error (MSE) loss to optimize quality prediction accuracy and the fidelity loss [39] to optimize quality monotonicity.

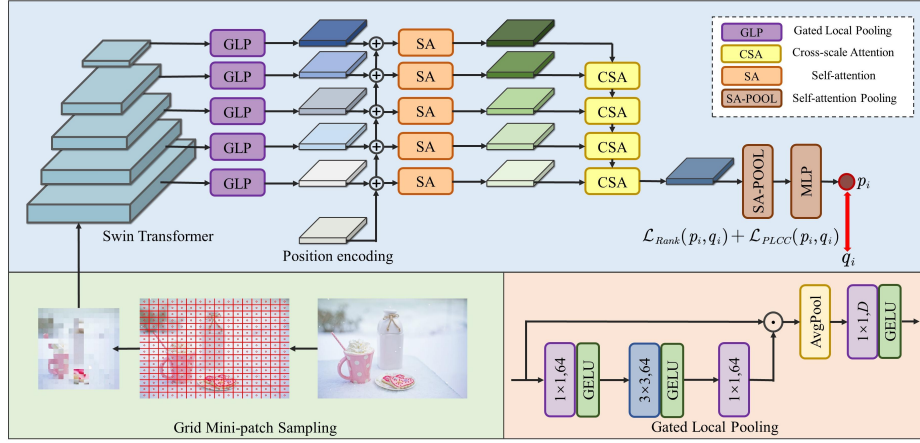
### 2.3 Blind Photo Quality Assessment based on Grid Mini-patch Sampling and Pyramid Perception

Songbai Tan<sup>1</sup>, Lixin Zhang<sup>2</sup>, Guanghui Yue<sup>2</sup>

<sup>1</sup> School of Management, Shenzhen University, China

<sup>2</sup> School of Biomedical Engineering, Shenzhen University, China  
Team SZU

We propose an effective photo quality assessment method named GS-PIQA, which is an improvement based on CFA-Net [5]. The detailed framework of the model is shown in Fig. 5. To enhance the ability to extract global information, we employ the Swin Transformer base network pre-trained on ImageNet as the backbone for GS-PIQA. In addition, GS-PIQA inherits the gated local pooling (GLP), the self-attention (SA) blocks, and the cross-scale attention (CSA) blocks in CFA-Net to enhance the multi-scale features across different layers. Through this top-down feature extraction and enhancement method, the model can form a pyramid perception capability. Given the high resolution of photos, directly resizing them would result in a significant loss of quality-related information. The common approach is to perform multiple crops on the image



**Fig. 5:** Overview of the proposed GS-PIQA by Team SZU.

and predict the quality for different cropped regions, averaging these regional quality scores to obtain the overall quality. While this method avoids the distortion and information loss caused by resizing, the small cropped areas can only represent local information, leading to substantial bias in overall quality prediction. To address these issues, we adopt the grid mini-patch sampling method for high-resolution images, which reduces the input resolution while preserving the semantic and quality features of the original image. Specifically, we cut the input high-resolution image  $\mathcal{P}$  into a uniform grid of  $N \times N$ , representing them as  $G = \{g_{(0,0)}, \dots, g_{(i,j)}, \dots, g_{(N-1,N-1)}\}$ , where  $i$  and  $j$  indicate that the grid is in the  $i$ -th row and  $j$ -th column, respectively. For each grid  $g_{(i,j)}$ , we randomly take a small region of size  $n \times n$  and splice all the obtained small regions to obtain the final sample image of size  $K \times K$ . In this experiment, the values of  $N$ ,  $n$ , and  $K$  are set to 16, 24, and 384, respectively. The uniform grid mini-patch sampling process is formalized as follows:

$$g_{(i,j)} = \mathcal{P}\left[\frac{i \times H}{N} : \frac{(i+1) \times H}{N}, \frac{j \times W}{N} : \frac{(j+1) \times W}{N}\right], \quad (1)$$

where  $H$  and  $W$  are the height and width of the input image respectively. The detail information of GS-PIQA is illustrated in Table. 4.

During the training process, we randomly sampled an image 10 times. The average of the quality prediction results of 10 samples is taken as the final quality prediction result of the image. To train the network, we employed the Rank and PLCC loss functions, which can be expressed as follows:

$$\mathcal{L} = \mathcal{L}_{\text{Rank}}(p_i, q_i) + \mathcal{L}_{\text{PLCC}}(p_i, q_i) \quad (2)$$

where  $p_i$  and  $q_i$  represent the predicted and true scores, respectively. Since the predicted results are not in the same range as the true quality scores, we map the predicted results as follows:

$$p_i = \frac{p_i - \min(p_i)}{\max(p_i) - \min(p_i)} \times (\max(q_i) - \min(q_i)) + \min(q_i) \quad (3)$$

**Implementation details** We trained and tested only on the UHD-IQA database and divided the database into training and test sets according to 8:2. The input images were processed using the grid mini-patch sampling method to obtain samples of size  $384 \times 384$ . To train GS-PIQA, we used the AdamW optimizer, initializing the learning rate at  $10^{-4}$  and the weight decay coefficient at  $10^{-5}$ . The network was trained for 10 epochs with the cosine learning rate decay strategy, setting the temperature coefficient  $T$  to 5.

The training process was divided into two phases. The first phase was trained using the above configuration, saving the results that performed best in the test set. In the second phase, we loaded the weights from the first phase and only fine-tuned the last fully connected layer. We increased the number of samples per image in the training set to 30. The fine-tuning learning rate was set to  $5 \times 10^{-5}$ , with a weight decay of  $10^{-5}$ , for training 10 epochs.

#### 2.4 High Resolution Patch Based Transformer with Quality-aware Feature Extraction

*Daekyu Kwon, Dongyoung Kim, Seon Joo Kim*  
 CIPLAB, Yonsei University, Korea

We propose a Vision Transformer [12] based IQA method which can efficiently handle arbitrary high-resolution images, using high-resolution patch strategy and quality-aware CNN extractor [20]. When applying the conventional ViT architecture to UHD images, excessive computation is required due to the large number of patches needed for training. To address this issue, we propose an architecture that can efficiently compute with fewer patches for high-resolution images by increasing the patch size, typically around 12 or 14, to 224. Doing so enables us to effectively handle UHD images with Vision Transformer architecture less than 50G MACs.

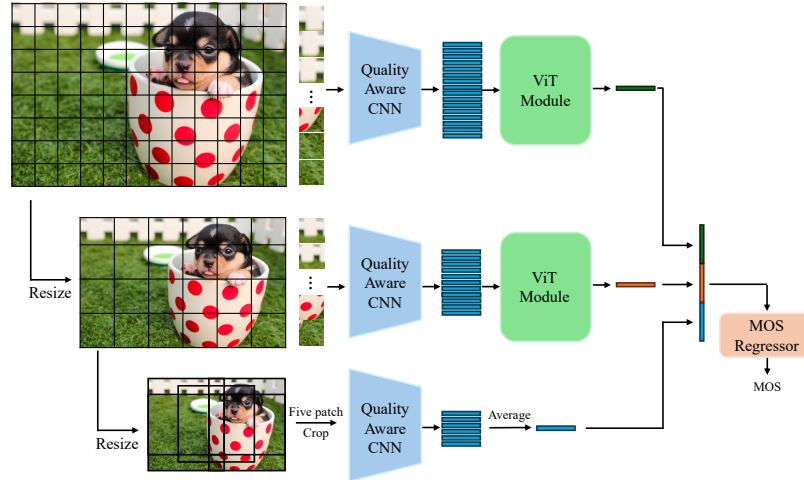
Furthermore, by employing high-resolution patches, we integrate an advanced CNN that can extract more meaningful features for IQA in the patch projection stage with ViT rather than the simple CNN utilized by the conventional ViT architecture. We first train a CNN-based feature extractor through a quality-aware pre-training method and utilize it as a feature extractor at the fine-tuning stage.

**Global Method Description** Our method consists of two primary stages: a pre-training stage (where we only train a CNN-based feature extractor) and a fine-tuning stage.

Our model consists of two primary components: a CNN-based feature extractor and a transformer-based feature aggregation module. For the CNN-based

Method	SRCC	PLCC
MobileNet(ImageNet-21k) + ViT	0.7828	0.7860
MobileNet(ATTIQA) + ViT	0.8063	0.8136

**Table 5:** Comparison of CIPLAB ensemble results using Quality-Aware CNN Extractor. We measure SRCC and PLCC using the official validation set.



**Fig. 6:** The overall process of the CIPLAB ensemble method. We utilize three types of various sized images. First, we patchify each image and encode them into features using a pre-trained quality-aware CNN. The features extracted from high-resolution images are encoded by a ViT module, while the features extracted from low-resolution images are averaged. We then concatenate these features and predict the MOS using a 2-layer MLP regressor.

feature extractor, we employ MobileNet-v3-large [17] from the timm library as the backbone and attach 2-layer MLPs to each attribution head, following the ATTIQA approach. For the transformer-based feature aggregation module, we utilize the default Vision Transformer architecture as a backbone, adopting Global Average Pooling for the final feature extraction instead of the CLS token. As we mentioned, we use a patch size of 224 and encode each patch into features using the CNN-based feature extractor.

**Pre-training Stage.** Our pre-training method is derived from ATTIQA [20]. Due to computational restrictions, we train MobileNet-V3 as a lightweight backbone using ATTIQA’s pre-training strategy with ImageNet-21k. We note that all pre-training setups are identical to ATTIQA’s setup, and additional details are provided in Section 3.

**Fine-tuning Stage.** Inspired by MUSIQ [19], we also utilize a multi-scale input strategy. To implement this strategy, we use three types of inputs: (a) the

original resolution image ( $W=3840$ ), (b) a 1/4 resolution image ( $W=960$ ), and (c) a tiny resolution image ( $W=256$ ). Each image is encoded into features independently using different CNN-based feature extractors.

Given that high-resolution images ((a) and (b)) are sufficient to integrate with the transformer, we extract features of high-resolution images using the transformer with images (a) and (b). For image (c), we compute the final feature by extracting five features for each side and center crop and averaging them. After extracting three features for high-resolution images and low-resolution image, we concatenate them into one feature and predict ground truth MOS using 2-layer MLP.

Details	Pre-train Stage	Finetune Stage
<b>Backbone</b>	MobileNetV3	MobileNetV3 + Vision Transformer
<b>Loss</b>	MarginRankingLoss	L1 Loss
<b>Optimizer</b>	AdamW	AdamW
<b>Learning Rate</b>	1e-4	1e-5
<b>GPU</b>	8 × V100	4 × RTX2080 Ti
<b>Dataset</b>	Imagenet 21k	UHD-IQA Dataset
<b>Times</b>	4d	10h
<b>Augmentation</b>	RandomResizedCrop	RandomHorizontalFlip(p=0.5)

**Table 6:** Implementation details for Pre-train and Finetune Stages of CIPLAB.

## 2.5 Learning from Strong to Weak, Enhanced Quality Comparison Network via Efficient Transfer Learning

*Yunchen Zhang, Xiangkai Xu, Hong Gao, Yiming Bao, Ji Shi, Xiugang Dong, Xiangsheng Zhou, Yaofeng Tu*  
ZTE Corporation

We propose two IQA models with different parameter scales. The teacher model, called Ensemble IQANet (EIQANet), is a large-parameter model designed to explore the upper bound of performance on UHD datasets [14]. The student model, Enhanced QCNet (EQCNet), is based on geometric order learning [21] for accurate rank estimation, serving as a lightweight model to meet the requirements of real-time applications. It is worth noticing that a significant performance gap lies in EIQANet and EQCNet. Furthermore, we designed a multi-stage knowledge transfer strategy involving three training steps: pre-training, fine-tuning, and calibration. This approach facilitates effective knowledge transfer between heterogeneous models and drives the construction of a well-arranged, well-clustered embedding space.

Method	KRCC	SROCC	RMSE	MAE
Q-Align [44]	0.3069	0.4412	0.1685	0.0289
Q-Align-LoRA (finetuned) [44]	0.2052	0.2624	0.0748	0.0597
Compare2Score [46]	0.2553	0.3735	0.1651	0.1524
QCN [32]	0.2977	0.42756	0.0615	0.0496
QCN-UHD (finetuned)	0.4707	0.6485	0.0581	0.0484
EQCN (Ours)	<b>0.6520</b>	<b>0.8403</b>	<b>0.0371</b>	<b>0.0289</b>

**Table 7:** Performance Comparisons of EQCN and latest BIQA methods.

**Enhanced IQANet** Inspired by RD-VQA [36], we propose the Enhanced IQANet (EIQANet). Given that the image resolutions in the UHD dataset all exceed 2K, we have made several advancements in image processing and feature extraction to fully utilize the information in high-resolution images. To better focus on the objective evaluation metrics of IQA tasks, we have also refined the loss functions. Our approach includes the following improvements:

**High-Resolution Image Processing.** The latest VLM model [9] processes images up to 4K resolutions without altering the image feature encoder architecture. We introduce the dynamic patch-slicing mechanism that allows a high-resolution image to be divided into up to 4 patches, capturing high-resolution features.

**Multi-Model Feature Fusion.** To boost the performance of the BIQA network, we introduce several advanced IQA models to provide auxiliary features:

QCN [32]: As the first image quality prediction model based on geometric order learning [21], QCN extracts features with strong generalization performance.

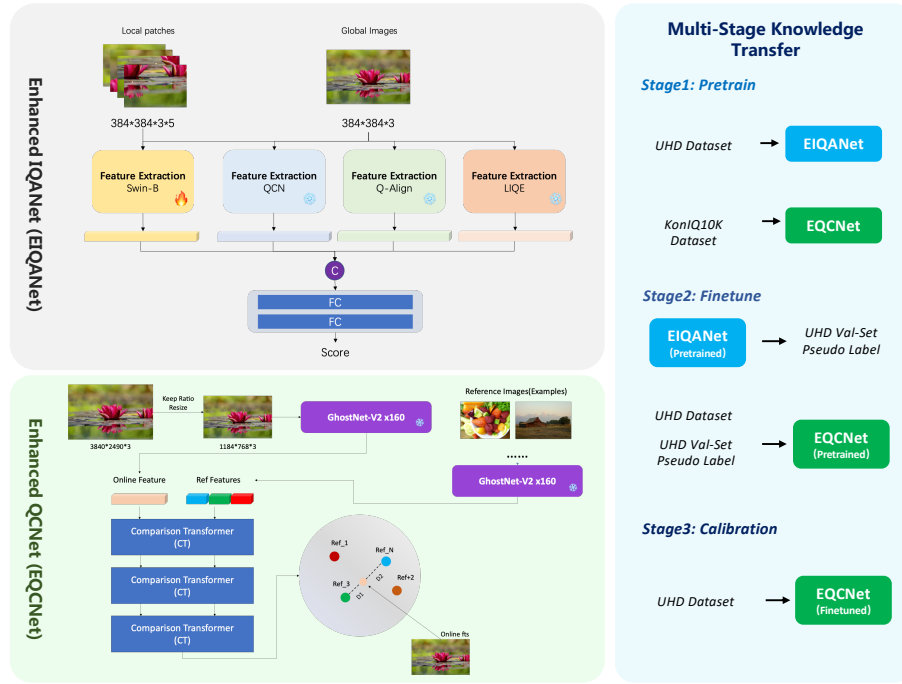
Q-Align [44]: As a VLM model, Q-Align leverages powerful LLMs to offer highly interpretable image quality assessments. We utilize the penultimate layer embeddings as features.

LIQE [45]: Based on image-text contrastive learning, the LIQE image encoder provides rich image features aligned with natural language.

Similar to RD-VQA, we employ an offline feature extraction method to obtain the above auxiliary features.

**Refined Loss Function.** [32] considered only the  $l1$  loss function during training, which overlooked the ordered sequence relationship of image quality within a batch, resulting in sub-optimal performance on objective evaluation metrics like PLCC. To address this, we additionally incorporate PLCC and SRCC loss functions, enabling the network to consider the absolute scores of the current samples and the relative order of image quality assessments within a batch.

**Enhanced QCNet (EQCNet)** The proposed EIQANet significantly improves performance metrics on the UHD dataset [14]. However, EIQANet’s reliance on offline feature extraction and its large number of parameters severely limit its practicality in real-world scenarios.



**Fig. 7:** Illustration of Enhanced IQANet (EIQANet), Enhanced Quality Comparison Network (EQCNet) and Multi-Stage Knowledge Transfer Strategy.

To address these limitations, following the design of [32], we introduce the comparison transformer (CT) to map each instance into a feature vector in an embedding space. Furthermore, the geometric order learning (GOL) [21] uses the reference points to satisfy both order and metric constraints and construct a well-arranged embedding space.

**Efficient Backbone Design.** To ensure computational efficiency when processing high-resolution images, we use the GhostNetV2 [38] as the backbone for image feature extraction. GhostNetV2, benefiting from the DFC attention mechanism [38] and depth-wise separable convolutions, ensures both feature diversity and model efficiency. We believe that GhostNetV2’s efficiency in modeling image features ensures that even after the features are projected through GOL, they retain their discriminative power.

**Multi-Stage Knowledge Transfer** Based on [21], the GOL method is significantly influenced by the initialization of reference points, which depend on the distribution characteristics of the provided training dataset. However, a substantial distribution difference exists between the original UHD training set and the test set data [14].

To address this, we designed a multi-stage knowledge transfer method. First, we pre-trained EQCNet using the KonIQ-10k [16] dataset to impart an initial image quality perception capability. Second, we utilized EIQANet, as mentioned in Sec. 2.5, to generate pseudo-labeled data on the validation set of the UHD dataset [14]. This pseudo-labeled data was then combined with the UHD training set for joint fine-tuning. Third, we fine-tuned EQCNet to align its embedding space with the joint UHD dataset distribution. Notably, EQCNet was initialized with weights from the model pre-trained on the KonIQ-10k dataset. Finally, recognizing potential noise and errors in the pseudo-labels, we further calibrated the EQCNet model using the UHD training set to obtain the final model.

This training method mitigates the slow convergence issue of small-parameter models and transfers knowledge from large-parameter models through a progressive learning strategy. This approach guides the EQCNet in learning a comprehensive feature mapping space, enhancing the performance and robustness of the BIQA model.

**Additional Implementation details** We implemented EIQANet and EQCNet using PyTorch. For EIQANet, we used the Adam optimizer with a learning rate of  $10^{-5}$  during the pre-training stage. Additionally, we trained the model 10 times and averaged the results to achieve robust score predictions.

The training strategy for EQCNet is more complex. During the pre-training stage, we used the AdamW optimizer with a learning rate of  $5 \times 10^{-5}$  and trained the model for 100 epochs on the KonIQ-10k dataset. For the fine-tuning and calibration stages, we switched to the Lion optimizer, setting the learning rates to  $3 \times 10^{-5}$  and  $5 \times 10^{-5}$ , respectively. The fine-tuning stage consisted of 100 epochs on the mixed dataset, while the calibration stage was limited to 20 epochs on the UHD train set.

## 2.6 MobileIQA: No-Reference Image Quality Assessment for Mobile Devices using Teacher-Student Learning

Zewen Chen<sup>1,2</sup>, Shunhan Xu<sup>3</sup>, Haochen Guo<sup>4</sup>, Yun Zeng<sup>5</sup>, Shuai Liu<sup>3</sup>, Jian Guo<sup>6</sup>, Juan Wang<sup>1</sup>, Bing Li<sup>1</sup>, Dehua Liu<sup>7</sup> and Hesong Liu<sup>7</sup>

<sup>1</sup> State Key Laboratory of Multimodal Artificial Intelligence Systems, CASIA

<sup>2</sup> School of Artificial Intelligence, University of Chinese Academy of Sciences

<sup>3</sup> College of Smart City, Beijing Union University

<sup>4</sup> College of Information and Electrical Engineering, Hebei University

<sup>5</sup> School of Economics and Management, China University of Petroleum-Beijing

<sup>6</sup> College of Robotics, Beijing Union University

<sup>7</sup> SHANGHAI TRANSSION INFORMATION TECHNOLOGY LIMITED

To address the challenge of high-resolution image quality assessment, we explore a structure based on MobileViT [24] and MobileNet [18] as backbone networks, namely MobileViT-IQA and MobileNet-IQA [8]. Inspired by the multiple



scores given by human annotators, we designed a multi-view opinion (MVO) module. This module can fuse the features extracted by the backbone network, simulating the assessment opinions of different annotators, and ultimately integrate them into an image quality score.

When dealing with high-resolution images, two challenges arise: (1) MobileViT demonstrates excellent performance but has high MACs, making it difficult to deploy on mobile devices; (2) MobileNet offers high computational efficiency, but its performance is not as robust as MobileViT. To address these issues, we employ knowledge distillation [7]. We first train a high-performance MobileViT-IQA model and then use it as a teacher model to guide the learning of the MobileNet-IQA. This model supports outputs with resolutions up to  $1907 \times 1231$  and requires only about 49 GMACs.

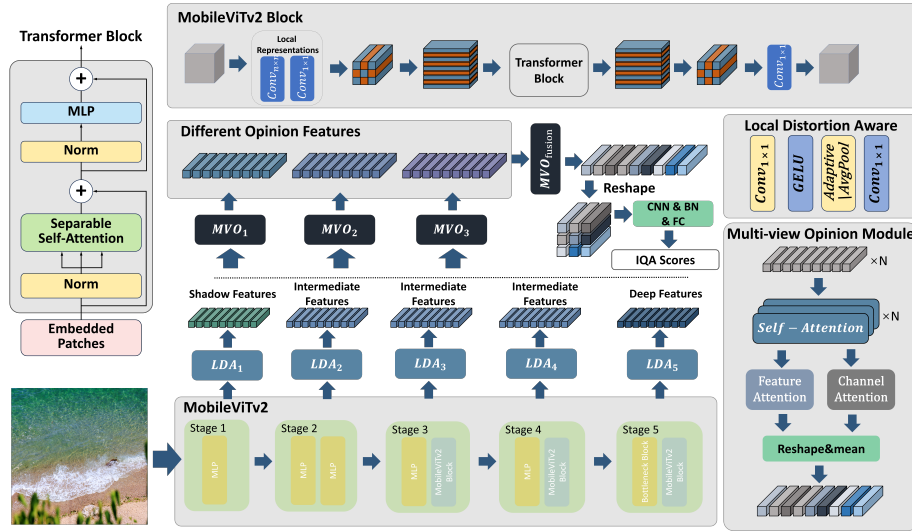
This approach effectively balances high performance and computational efficiency, providing a viable solution for high-resolution image quality assessment on mobile devices.

**Model Design** We take the features captured from five layers in the MobileViT and MobileNet. Many existing works prove that the multi-layer features are helpful for the IQA task [6, 7, 35, 41].

The teacher model (MobileViT-IQA) is shown in Fig 8. First, multi-scale features are extracted from five layers of MobileViT, enabling the model to comprehend image quality more comprehensively. Subsequently, these features are fused and dimensionally reduced through a Local Distortion Aware (LDA) module. The processed five features are then input into three Multi-view Opinion (MVO) modules with different weight initializations, generating three distinct opinion features that simulate subjective opinions of the same image by multiple assessors. Finally, these three opinion features are integrated through an additional MVO module, followed by reshaping, convolutional neural network (CNN), and fully connected (FC) layer operations to derive the final image quality score. The student model (MobileNet-IQA) shares the same framework as MobileViT-IQA but uses MobileNet as the backbone.

The distillation process is shown in Fig 9. Since the MobileViT-IQA and the MobileNet-IQA share the same framework, distilling the teacher’s knowledge to the student is more efficient. We take the *MSE* loss to supervise the discrepancy between the Different Opinion Features (DOF) from the teacher and student models. During training, the discrepancy between the predicted and GT scores is also supervised by the *MSE* loss.

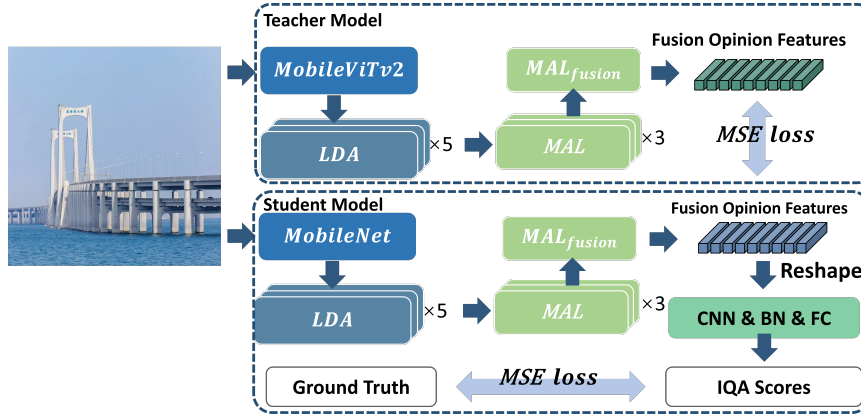
**Multi-view Opinion** The motivation is that individuals often have diverse subjective perceptions and regions of interest when viewing the same image. To this end, we employ multiple MVOs to learn attention from different viewpoints. Each MVO is initialized with different weights and updated independently to encourage diversity and avoid redundant output features. The number of MVOs can be flexibly set as a hyper-parameter. In this work, we set the number to 3. As shown in Fig 8, the MVO starts from  $N$  self-attentions (SAs), each of



**Fig. 8:** Framework of the teacher model MobileViT-IQA [8]. The student model MobileNet-IQA shares the same framework, with MobileNet as its backbone.

which is responsible for processing a basic feature  $\mathbf{f}_j$  ( $1 \leq j \leq N$ ). The outputs of all the SAs are concatenated, forming a multi-level aggregated feature  $\mathbf{F} \in \mathbb{R}^{C \times D \times N}$ . Then  $\mathbf{F}$  passes through two branches, i.e., a pixel-wise SA branch and a channel-wise SA branch, which apply an SA across spatial and channel dimensions, respectively, to capture complementary non-local contexts and generate multi-view attention maps. In particular, for the channel-wise SA, the feature  $\mathbf{F}$  is first reshaped and permuted to convert the size from  $C \times D \times N$  to  $D \times (C \times N)$ . After the SA, the output feature is permuted and reshaped back to the original size  $C \times D \times N$ . Subsequently, the outputs of the two branches are added and average pooled, generating an opinion feature. The design of the two branches has two key advantages. First, implementing the SA in different dimensions promotes diverse attention learning, yielding complementary information. Second, contextualized long-range relationships are aggregated, benefiting global quality perception.

**Image Quality Score Regression.** Assuming that  $M$  opinion features are generated from  $M$  MVOs. To derive a global quality score from the collected opinion features, we utilize an additional MVO. The MVO integrates diverse contextual perspectives, resulting in a comprehensive opinion feature that captures essential information. This feature is then processed through a transformer block, three convolutional layers with kernel sizes of  $5 \times 5$ ,  $3 \times 3$ , and  $3 \times 3$  to reduce the number of channels, followed by two fully connected layers that transform the feature size from 128 to 64 and from 64 to 1. Finally, we obtain a predicted quality score.



**Fig. 9:** Model distillation process of MobileNet-IQA. The purpose of teacher-student learning is achieved by supervising the Different Opinion Features of the teacher and student networks.

**Additional Implementation Details** We use the MSE loss to reduce the discrepancy between predicted and GT scores. Then, we use the Adam optimizer with a learning rate of  $10^{-5}$  and a weight decay of  $10^{-5}$ . The learning rate is adjusted using the Cosine Annealing for every 50 epochs. We train the teacher model for 100 epochs (about 18h) with a batch size of 4 and the student model for 300 epochs (about 48h) with a batch size of 8 on one NVIDIA RTX800.

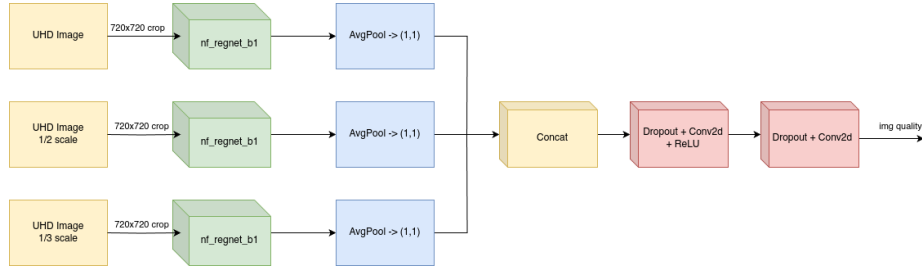
## 2.7 Multi-scale NF-RegNets Ensemble

*Grigory Malivenko*

The solution contains three parts, as well as the fusion block. Each sub-model is a NFRegNet [4] (Norm-Free RegNet) model (*nf-regnet-b1*) trained to predict photo quality on a specific resolution (1:1, 1:2, and 1:3). Features of these models are being fused together and used for the final photo quality estimation.

Without TTAs (test-time augmentations), it takes only 19.08 GMACs to process a photo. Each sub-model takes around 6.36 GMACs to run, and the fusion/classification block takes 0.74 MMACs. This fact makes it possible to perform TTAs very effectively: calculate features for all sub-models separately and then use the fusion/classification block for each possible combination. The runtime is 40 ms for each photo with TTAs, and 15 ms without TTAs.

*Implementation details* PyTorch with Adam optimizer was used. Standard learning rate  $10^{-3}$  with step lr-scheduler was used (every 15 steps, factor 0.8). Every sub-model was trained for 150 epochs. Then, after merging sub-models into a single model, only the fusion block was trained for 20 epochs (while sub-model



**Fig. 10:** Multi-scale NF-RegNets ensemble solution.

weights were frozen). Finally, the whole model was fine-tuned for another 20 epochs.

Each sub-model training process took around 2 hours, and initial fusion block training took around 30 minutes. The whole model needed to be fine-tuned for another 2 hours. Only random crops and random flips were used for augmentation.

For the final version of the solution, the model was trained on a whole dataset and fine-tuned on a pseudo-labeled validation part.

## 2.8 Hybrid Local-Global Image Quality Assessment

*Xingyuan Ma, Cheng Li*

We divide the original image into several patches and score them separately. To avoid the impact of image content on model performance, we randomly disrupt the order of the above patches and reassemble new images for scoring. Finally, the scores of the original image, several patches, and the reorganized new image are averaged to create the final score.

**Global Method Description** Our method, denoted as CLIP-IQA\*, is based on CLIP-IQA. Unlike CLIP-IQA, we use positional encoding, and the model’s input is fixed to  $224 \times 224$  pixels.

The prompts we used are ‘The quality of this photo is bad’, ‘The quality of this photo is poor’, ‘The quality of this photo is fair’, ‘The quality of this photo is good’, ‘The quality of this photo is perfect’.

In this challenge, an image resolution that is too large will require considerable calculation. A common approach is to downsample the original image to a very small resolution, such as  $224 \times 224$ , resulting in a severe loss of input information. In addition, this method violates the logic of subjective image quality evaluation. Inspired by the process of image quality evaluation from the whole to the part or from the part to the whole, we divide the original image into several patches and score them separately. The dimensions of image quality evaluation

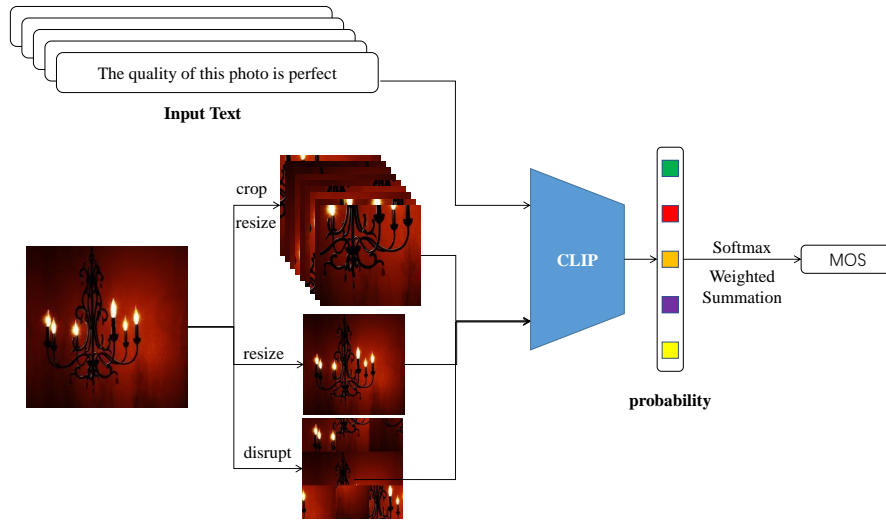


Fig. 11: The diagram of the proposed CLIP-IQA\*.

are related to noise, clarity, color, details, etc. To avoid the impact of image content on model performance, we randomly disrupt the order of the above patches and reassemble new images for scoring. Finally, the scores of the original image, several patches, and the reorganized new image are averaged as the final score.

**Implementation Details** During training and testing, the input data is processed as follows. First, we evenly divide the original image into 9 patches. Secondly, we shuffle the order of the 9 patches and reassemble them into a new image of the original image size. Then, we resize the original image, 9 patches, and the reorganized image to  $224 \times 224$ . Finally, all the above images are input into the model, and the scores of each image are averaged as the final score.

During training, the batch size is 3, and the total epochs are set to 80. We use Smooth-L1 loss as the training loss and CosineAnnealingLR for learning rate decay. In addition, the model with the highest MOS on the validation set is finally selected for testing.

## 2.9 Blind IQA Using Multiple Vision Encoders

Joonhee Lee<sup>1</sup>, Junseo Bang<sup>1</sup>, Se Young Chun<sup>1,2</sup>

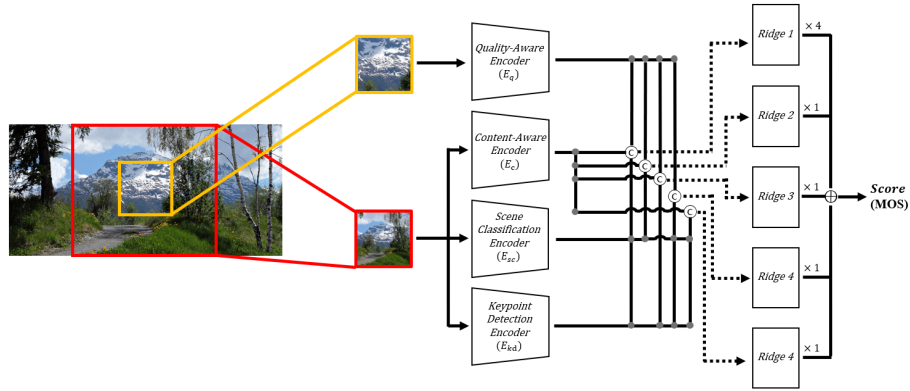
<sup>1</sup> Department of Electrical and Computer Engineering,

<sup>2</sup> INMC, Interdisciplinary Program in AI,

Seoul National University, Republic of Korea

Team ICL

In this study, we demonstrate that utilizing various image representations enhances perceptual understanding of images and improves the prediction of



**Fig. 12:** ICL team overall architecture. We use four pre-trained encoders as feature extractors and five ridge regressors to map these features to quality predictions.

Image Quality Assessment (IQA) scores. Four pre-trained encoders are employed as feature extractors, and five Ridge regressors are used to map these features to quality predictions. Specifically, along with the Quality-Aware Encoder and Content-Aware Encoders derived from the existing Re-IQA [31], we added task-specific encoders beneficial to IQA. Using this concept, we calculated the IQA score by linearly summing the outputs from the regressors.

The training dataset consisted solely of the 4K images provided by the challenge. However, using images of large size as input exceeded the computational limits set by the challenge. Therefore, a pre-processing step was implemented to crop the center of images to  $320 \times 320$  pixels before feeding them into the model. The cropping method varied depending on the encoder requirements. For encoders that required global information (content-aware, scene classification, keypoint detection), the images were first cropped to the largest possible square and then resized to  $320 \times 320$  pixels. For encoders that required local information (quality-aware), patches of  $320 \times 320$  pixels were employed without resizing.

During training, the features of the four encoders were regressed using ridge regressors. Five ridge regressors were trained; one regressed the features from all encoders, while the other four regressed combinations of features from three encoders each. During inference, the features from the four pre-trained encoders were passed through the five ridge regressors to yield five scores. Each score was weighted and combined to determine the final score (MOS).

**Implementation Details** We optimized the Ridge regression model using the “GridSearchCV” function of Scikit-learn [30]. The hyperparameter alpha was scanned from  $10^{-6}$  to  $10^6$ , with 13 equally spaced values on a log scale. We used the entire challenge dataset, dividing the labeled training dataset into a

0.8/0.2 split for training and validation data. With only the ridge regressors being optimized, the training time took approximately 5 to 6 minutes based on NVIDIA A100, the number of parameters is 139.1M, and MACs are 42.09G.

## Acknowledgements

This work was partially supported by the Humboldt Foundation. We thank the AIM 2024 sponsors: Meta Reality Labs, KuaiShou, Huawei, Sony Interactive Entertainment, and the University of Würzburg (Computer Vision Lab).

## References

1. Agnolucci, L., Galteri, L., Bertini, M.: Quality-Aware Image-Text Alignment for Real-World Image Quality Assessment. arXiv preprint arXiv:2403.11176 (2024) [5](#), [6](#)
2. Agnolucci, L., Galteri, L., Bertini, M., Del Bimbo, A.: ARNIQA: Learning Distortion Manifold for Image Quality Assessment. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 189–198 (2024) [5](#), [6](#)
3. Borji, A., Itti, L.: State-of-the-art in Visual Attention Modeling. IEEE Transactions on Pattern Analysis and Machine Intelligence **35**(1), 185–207 (2012) [9](#)
4. Brock, A., De, S., Smith, S.L., Simonyan, K.: High-Performance Large-Scale Image Recognition Without Normalization (2021), <https://arxiv.org/abs/2102.06171> [19](#)
5. Chen, C., Mo, J., Hou, J., Wu, H., Liao, L., Sun, W., Yan, Q., Lin, W.: TOPIQ: A Top-down Approach from Semantics to Distortions for Image Quality Assessment. IEEE Transactions on Image Processing (2024) [9](#)
6. Chen, Z., Qin, H., Wang, J., Yuan, C., Li, B., Hu, W., Wang, L.: PromptIQA: Boosting the Performance and Generalization for No-Reference Image Quality Assessment via Prompts. arXiv Preprint arXiv:2403.04993 (2024) [17](#)
7. Chen, Z., Wang, J., Li, B., Yuan, C., Xiong, W., Cheng, R., Hu, W.: Teacher-Guided Learning for Blind Image Quality Assessment. In: Proceedings of the Asian Conference on Computer Vision. pp. 2457–2474 (2022) [17](#)
8. Chen, Z., Xu, S., Zeng, Y., Guo, H., Guo, J., Liu, S., Wang, J., Li, B., Hu, W., Liu, D., et al.: Mobileiqa: Exploiting mobile-level diverse opinion network for no-reference image quality assessment using knowledge distillation. arXiv preprint arXiv:2409.01212 (2024) [16](#), [18](#)
9. Chen, Z., Wang, W., Tian, H., Ye, S., Gao, Z., Cui, E., Tong, W., Hu, K., Luo, J., Ma, Z., et al.: How Far Are we to GPT-4V? Closing the Gap to Commercial Multimodal Models with Open-source Suites. arXiv preprint arXiv:2404.16821 (2024) [14](#)
10. Conde, M.V., Lei, Z., Li, W., Bampis, C., Katsavounidis, I., Timofte, R., et al.: AIM 2024 Challenge on Efficient Video Super-Resolution for AV1 Compressed Content. In: Proceedings of the European Conference on Computer Vision (ECCV) Workshops (2024) [4](#)
11. Conde, M.V., Vasluianu, F.A., Xiong, J., Ye, W., Ranjan, R., Timofte, R., et al.: Compressed Depth Map Super-Resolution and Restoration: AIM 2024 Challenge Results. In: Proceedings of the European Conference on Computer Vision (ECCV) Workshops (2024) [4](#)

12. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Deghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In: ICLR (2021) [11](#)
13. Gu, J., Cai, H., Dong, C., Ren, J.S., Timofte, R., Gong, Y., Lao, S., Shi, S., Wang, J., Yang, S., et al.: Ntire 2022 challenge on perceptual image quality assessment. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 951–967 (2022) [4](#)
14. Hosu, V., Agnolucci, L., Wiedemann, O., Iso, D.: UHD-IQA Benchmark Database: Pushing the Boundaries of Blind Photo Quality Assessment. arXiv Preprint arXiv:2406.17472 (2024) [1](#), [13](#), [14](#), [15](#), [16](#)
15. Hosu, V., Conde, M.V., Agnolucci, L., Barman, N., Zadtootaghaj, S., Timofte, R., et al.: AIM 2024 challenge on uhd blind photo quality assessment. In: Proceedings of the European Conference on Computer Vision (ECCV) Workshops (2024) [4](#)
16. Hosu, V., Lin, H., Sziranyi, T., Saupe, D.: KonIQ-10k: An Ecologically Valid Database for Deep Learning of Blind Image Quality Assessment. IEEE Transactions on Image Processing **29**, 4041–4056 (2020) [16](#)
17. Howard, A., Sandler, M., Chu, G., Chen, L.C., Chen, B., Tan, M., Wang, W., Zhu, Y., Pang, R., Vasudevan, V., et al.: Searching for MobileNetV3. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1314–1324 (2019) [4](#), [12](#)
18. Howard, A.G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., Adam, H.: MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. arXiv Preprint arXiv:1704.04861 (2017) [4](#), [16](#)
19. Ke, J., Wang, Q., Wang, Y., Milanfar, P., Yang, F.: MUSIQ: Multi-scale Image Quality Transformer. In: ICCV. pp. 5148–5157 (2021) [12](#)
20. Kwon, D., Kim, D., Ki, S., Jo, Y., Lee, H.E., Kim, S.J.: CLIP-Guided Attribute Aware Pretraining for Generalizable Image Quality Assessment. arXiv preprint arXiv:2406.01020 (2024) [4](#), [11](#), [12](#)
21. Lee, S.H., Shin, N.H., Kim, C.S.: Geometric Order Learning for Rank Estimation. Advances in Neural Information Processing Systems **35**, 27–39 (2022) [13](#), [14](#), [15](#)
22. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10012–10022 (2021) [8](#), [9](#)
23. Madhusudana, P.C., Birkbeck, N., Wang, Y., Adsumilli, B., Bovik, A.C.: Image Quality Assessment Using Contrastive Learning. IEEE Transactions on Image Processing **31**, 4149–4161 (2022) [5](#), [6](#)
24. Mehta, S., Rastegari, M.: MobileVit: Light-Weight, General-Purpose, and Mobile-Friendly Vision Transformer. arXiv Preprint arXiv:2110.02178 (2021) [4](#), [16](#)
25. Molodetskikh, I., Borisov, A., Vatolin, D.S., Timofte, R., et al.: AIM 2024 Challenge on Video Super-Resolution Quality Assessment: Methods and Results. In: Proceedings of the European Conference on Computer Vision (ECCV) Workshops (2024) [4](#)
26. Moskalenko, A., Bryntsev, A., Vatolin, D.S., Timofte, R., et al.: AIM 2024 challenge on video saliency prediction: Methods and results. In: Proceedings of the European Conference on Computer Vision (ECCV) Workshops (2024) [4](#)
27. Murray, N., Marchesotti, L., Perronnin, F.: AVA: A Large-scale Database for Aesthetic Visual Analysis. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 2408–2415 (2012) [8](#)



28. Nazarczuk, M., Catley-Chandar, S., Tanay, T., Shaw, R., Pérez-Pellitero, E., Timofte, R., et al.: AIM 2024 Sparse Neural Rendering Challenge: Methods and Results. In: Proceedings of the European Conference on Computer Vision (ECCV) Workshops (2024) [4](#)
29. Nazarczuk, M., Tanay, T., Catley-Chandar, S., Shaw, R., Timofte, R., Pérez-Pellitero, E.: AIM 2024 Sparse Neural Rendering Challenge: Dataset and Benchmark. In: Proceedings of the European Conference on Computer Vision (ECCV) Workshops (2024) [4](#)
30. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al.: Scikit-learn: Machine Learning in Python. *The Journal of Machine Learning Research* **12**, 2825–2830 (2011) [22](#)
31. Saha, A., Mishra, S., Bovik, A.C.: Re-IQA: Unsupervised Learning for Image Quality Assessment in the Wild. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5846–5855 (2023) [22](#)
32. Shin, N.H., Lee, S.H., Kim, C.S.: Blind Image Quality Assessment Based on Geometric Order Learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2024) [14](#), [15](#)
33. Smirnov, M., Gushchin, A., Antsiferova, A., Vatolin, D.S., Timofte, R., et al.: AIM 2024 Challenge on Compressed Video Quality Assessment: Methods and Results. In: Proceedings of the European Conference on Computer Vision (ECCV) Workshops (2024) [4](#)
34. Su, S., Yan, Q., Zhu, Y., Zhang, C., Ge, X., Sun, J., Zhang, Y.: 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 (CVPR). pp. 3667–3676 (2020) [5](#), [6](#)
35. Su, S., Yan, Q., Zhu, Y., Zhang, C., Ge, X., Sun, J., Zhang, Y.: Blindly assess image quality in the wild guided by a self-adaptive hyper network. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2020) [17](#)
36. Sun, W., Wu, H., Zhang, Z., Jia, J., Zhang, Z., Cao, L., Chen, Q., Min, X., Lin, W., Zhai, G.: Enhancing Blind Video Quality Assessment with Rich Quality-aware Features. arXiv preprint arXiv:2405.08745 (2024) [14](#)
37. Sun, W., Zhang, W., Cao, Y., Cao, L., Jia, J., Chen, Z., Zhang, Z., Min, X., Zhai, G.: Assessing uhd image quality from aesthetics, distortions, and saliency. arXiv preprint arXiv:2409.00749 (2024) [8](#)
38. Tang, Y., Han, K., Guo, J., Xu, C., Xu, C., Wang, Y.: GhostNetv2: Enhance Cheap Operation with Long-Range Attention. *Advances in Neural Information Processing Systems* **35**, 9969–9982 (2022) [15](#)
39. Tsai, M.F., Liu, T.Y., Qin, T., Chen, H.H., Ma, W.Y.: Frank: a Ranking Method With Fidelity Loss. In: Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 383–390 (2007) [8](#), [9](#)
40. Wang, J., Chan, K.C., Loy, C.C.: Exploring CLIP for Assessing the Look and Feel of Images. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 37, pp. 2555–2563 (2023) [5](#), [6](#)
41. Wang, J., Chen, Z., Yuan, C., Li, B., Ma, W., Hu, W.: Hierarchical Curriculum Learning for No-Reference Image Quality Assessment. *International Journal of Computer Vision* **131**(11), 3074–3093 (2023) [17](#)
42. Wiedemann, O., Hosu, V., Su, S., Saupe, D.: KonX: Cross-Resolution Image Quality Assessment. *Quality and User Experience* **8**(1), 8 (Dec 2023). <https://doi.org/10.1007/s41233-023-00061-8> [5](#), [6](#)

43. Wu, H., Chen, C., Hou, J., Liao, L., Wang, A., Sun, W., Yan, Q., Lin, W.: Fast-VQA: Efficient End-to-end Video Quality Assessment with Fragment Sampling. In: European Conference on Computer Vision. pp. 538–554 (2022) [9](#)
44. Wu, H., Zhang, Z., Zhang, W., Chen, C., Liao, L., Li, C., Gao, Y., Wang, A., Zhang, E., Sun, W., et al.: Q-Align: Teaching LLMs for Visual Scoring via Discrete Text-defined Levels. arXiv preprint arXiv:2312.17090 (2023) [14](#)
45. Zhang, W., Zhai, G., Wei, Y., Yang, X., Ma, K.: Blind Image Quality Assessment via Vision-Language Correspondence: A Multitask Learning Perspective. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 14071–14081 (2023) [14](#)
46. Zhu, H., Wu, H., Li, Y., Zhang, Z., Chen, B., Zhu, L., Fang, Y., Zhai, G., Lin, W., Wang, S.: Adaptive Image Quality Assessment via Teaching Large Multimodal Model to Compare. arXiv preprint arXiv:2405.19298 (2024) [14](#)