

LAR-IQA: A Lightweight, Accurate, and Robust No-Reference Image Quality Assessment Model

Nasim Jamshidi Avanaki¹, Abhijay Ghildyal¹, Nabajeet Barman², and Saman Zadtootaghaj²

¹ Portland State University, OR, USA

² Sony Interactive Entertainment, UK/Germany

n.jamshidi.avanaki@gmail.com, abhijay@pdx.edu
{Nabajeet.Barman,Saman.Zadtootaghaj}@sony.com

Abstract. Recent advancements in the field of No-Reference Image Quality Assessment (NR-IQA) using deep learning techniques demonstrate high performance across multiple open-source datasets. However, such models are typically very large and complex making them not so suitable for real-world deployment, especially on resource- and battery-constrained mobile devices. To address this limitation, we propose a compact, lightweight NR-IQA model that achieves state-of-the-art (SOTA) performance on ECCV AIM UHD-IQA challenge validation and test datasets while being also nearly 5.7 times faster than the fastest SOTA model. Our model features a dual-branch architecture, with each branch separately trained on synthetically and authentically distorted images which enhances the model’s generalizability across different distortion types. To improve robustness under diverse real-world visual conditions, we additionally incorporate multiple color spaces during the training process. We also demonstrate the higher accuracy of recently proposed Kolmogorov-Arnold Networks (KANs) for final quality regression as compared to the conventional Multi-Layer Perceptrons (MLPs). Our evaluation considering various open-source datasets highlights the practical, high-accuracy, and robust performance of our proposed lightweight model. Code: <https://github.com/nasimjamshidi/LAR-IQA>.

Keywords: Quality Assessment · IQA · No-Reference IQA · BIQA · Real-time · Lightweight

1 Introduction

The Image Quality Assessment (IQA) task of measuring the quality of images as perceived by humans remains one of the most interesting and challenging fields in computer vision. No-Reference IQA (NR-IQA), also known as Blind IQA (BIQA), focuses on estimating the quality of degraded images when there is no high-quality reference image available for comparison. This task is particularly challenging given the wide range of possible distortions (compression, blur, noise, etc.) that might be present in an image. NR-IQA plays a critical role across multiple industries and applications, such as photography, surveillance, healthcare, automotive, social media, and user-generated content platforms. Millions

of user-generated content (UGC) images are uploaded daily and shared across numerous social media platforms such as Instagram, X, and Flickr. In the wild, user-captured images can suffer from distortions such as blurriness, noise (from the camera sensor), color distortions, compression artifacts (blockiness), or a combination of these issues. Automatically detecting low-quality or inappropriate images and guiding the necessary pre- and post-processing steps (quality enhancement, compression factor, deblurring, etc.) is critical for enhancing user experience and the success of such companies.

One of the major challenges in NR-IQA is developing smaller, faster and more efficient methods suitable for real-time quality assessment. Traditional NR-IQA models are generally fast and less complex but often lack accuracy [20, 27–30, 33, 44, 45, 47, 49]. On the other hand, traditional DNN-based models, while more accurate, typically have higher complexity and are computationally intensive [3, 5, 11, 35, 39, 46, 48]. Recent IQA models based on large multi-modal models utilize sophisticated architectures such as transformers for both vision and text encoders [43, 51] resulting in large model size and complexity, making them unsuitable for most real-time evaluation tasks.

Some more recent NR-IQA methods use Transformers as their backbone network [11, 46, 48] since Transformers have been shown to provide better features for NR-IQA, resulting in more accurate and robust results. However, Transformer-based models consists of a large number of parameters and require significant computational resources, both for training and inference, thus restricting their applicability for deployment on low-power devices. Recently, Vision Language Models (VLMs), particularly CLIP [32], have achieved significant success in NR-IQA. When fine-tuned for NR-IQA tasks [2, 12, 39], CLIP demonstrates good accuracy, robustness, and generalization compared to existing methods, effectively capturing a wide range of diverse distortions in large datasets. However, we do not utilize VLM-based pre-trained networks and instead focus on developing an accurate and robust model that is more computationally efficient, specifically in terms of the number of Multiply-Accumulate (MAC) operations required for a forward pass with an input image size of 3840×2160 pixels.

1.1 Contributions

This paper makes several contributions for low complexity, generalizability and robustness which is discussed next.

Lower Complexity. Due to the ability of DNN-based methods to capture intricate patterns and non-linear distortions for higher accuracy and robustness, we use the lightweight MobileNetV3 [16] as the baseline network. Compared to other DNN-based NR-IQA methods [2, 3, 35, 39], our network has fewer parameters, resulting in lower complexity.

Generalizability. One of the main challenges in training a generalizable NR-IQA model is that synthetic datasets lack the diverse distortions found in real-world scenarios. Several methods have attempted to address this gap using pre-training, self-supervision, and novel loss functions [2, 3, 11, 39]. Self-supervised learning with contrastive loss or pre-training on unlabeled data [3, 24] is effective

for handling both synthetic and authentic types of degradation. However, we address the problem of lack of comprehensive datasets by training two separate models — one on dataset with synthetic distortions and another on authentic distortion datasets. Both individually trained models are then combined to create a more generalized model.

Such multiple branch architectures have been proposed previously. For example, authors in [12] used three branches: semantic, aesthetic, and technical, the prediction scores from which are then fused together to obtain the final quality score. In contrast, ours is the first work to propose a dual-branch architecture to tackle two different types of distortions, synthetic and authentic. Given the need for real-time applications, we focus on making our NR-IQA model both lightweight and fast by focusing exclusively on MobileNet [16, 22, 38, 50] as the backbone image encoder for our model.

Robustness and Accuracy. To enhance our lightweight model’s robustness and accuracy, we incorporate several strategies: separate dual branch training (synthetic and authentic distortions), using KAN [21] as the regression head, and using a color space loss function.

In summary, this paper introduces a novel approach and model for low-complexity NR-IQA, leveraging both authentic and synthetic distortions for robust and accurate evaluations. Our key contributions are:

- We propose combining authentic and synthetic IQA branches, each trained on different datasets, to enhance the model’s robustness and generalizability across various distortions.
- We compare different backbone architectures tailored for the synthetic and authentic branches, facilitating optimal design choices.
- We conduct an ablation study comparing KANs and MLPs for the quality regression module. To the best of our knowledge, this is the first study to investigate and compare the efficacy of KANs against MLPs in this context.
- We propose a novel loss function addressing variations in different color spaces to improve the robustness of the IQA models.

Our approach results in a lightweight model that surpasses state-of-the-art performance [13] on the validation and test datasets of the ECCV AIM UHD-IQA challenge³ while maintaining efficiency for practical applications.

The rest of the paper is organized as follows: In Section 2, we discuss prior work related to No-Reference Image Quality Assessment and the integration of synthetic and authentic data. In Section 3, we provide a detailed description of our model architecture and the comparative analysis of different backbones, heads, and loss functions. Section 4 presents the results of our ablation study and the evaluation of our proposed model versus SOTA. Finally, Section 5 concludes this paper by highlighting the implications of our findings.

³ <https://codalab.lisn.upsaclay.fr/competitions/19335>

2 Related Work

2.1 No-Reference Image Quality Assessment

Over the years, many NR-IQA models have been proposed, from initial traditional approaches based on Natural Scene Statistics (NSS) and hand-crafted features to deep-learning based approaches to more recently, models based on Vision-Language Models (VLMs). Traditional NR-IQA approaches relied on the detection and measurement of specific distortions such as blockiness, blur, banding, and noise, among others. Such methods leveraged hand-crafted natural scene statistics (NSS) features to assess image quality [20, 27–30, 33, 40, 44, 45, 47, 49]. However, these distortion-specific models based on hand-crafted features do not generalize well to images with multiple types of distortions.

Recently, driven by the success of deep neural networks (DNNs) in computer vision, several NR-IQA methods have been proposed that generalize better without the need for explicitly designing features. DNN-based approaches achieve high accuracy across various datasets by learning complex patterns and accounting for multiple distortions in images [3, 5, 46, 48]. However, due to the limited size of available datasets, they often tend to overfit on training data and cannot generalize well to newer distortions (not present in the training data), which are evident when performing cross-dataset evaluation as discussed in [11].

More recently, Vision-Language Models (VLMs) have found success in the field of IQA, combining both visual and textual information to assess the image quality [39]. Since VLMs have been trained on very large-scale datasets that combine vision and language modalities, they have a more nuanced understanding of the image content using its context. A combination of semantic information and textual descriptions/captions allows them to better understand the quality, often resembling human ratings. Similarly, large multi-modality models (MLLM) have been developed for the IQA task [43, 51] and are more accurate than traditional DNN-based methods. However, their complex training setups and large scale make VLM and MLLM-based models computationally expensive.

2.2 Generalization Across Diverse Datasets

Training large Vision Transformer (ViT) or CNN backbone networks requires extensive datasets for pre-training to ensure generalization to new datasets. While some approaches utilize contrastive learning for pre-training [3, 24], others attempt to merge multiple image datasets with available subjective scores [26]. However, merging datasets is challenging due to inherent subjective biases within each dataset. As described in [26], these biases include rating noise, subjective test order effects, varying distributions of quality scores across datasets, and long-term dependencies in subjective experiments. Additional biases, especially when crowdsourcing is used, arise from less controlled environments, such as differences in monitor distances, display sizes, and settings (e.g., gamma and luminance variations). These biases complicate the straightforward merging of subjective scores from different tests.

To address these challenges, various alternative methods have been proposed, such as integrating subjective biases into the loss function [26], and using ranking

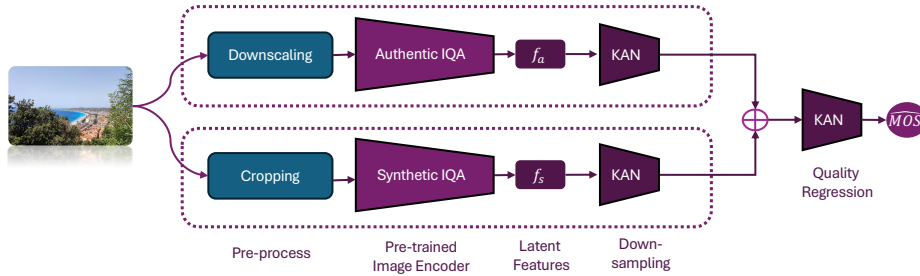


Fig. 1: Proposed model architecture. The image quality is evaluated using two branches: Authentic and Synthetic. In both branches, MobileNetV3 [16] serves as the lightweight image encoder. The features extracted from these branches are concatenated and then used as input to KAN, the quality regression module, which outputs the final predicted image quality score.

loss instead of Mean Squared Error (MSE) loss, based on the assumption that subjective biases do not affect the ranking of MOS within a dataset [11]. We on the other hand, train each branch of our model using multiple authentic and synthetic datasets by employing a multi-task training approach. For each branch, each dataset is considered as a separate task with a unique head. This allows each branch to effectively learn the subjective biases present in respective individual datasets leading to a better generalization of the image encoder. In the final model, we remove these individual heads, and instead use a single KAN head for each of the two branches as a down-sampler, and then fuse the resulting embeddings with a larger KAN head.

2.3 Kolmogorov-Arnold Networks (KAN)

While typically fully-connected MLPs have been used as regression heads [7, 12], recently, KANs [21] have been introduced as an alternative to MLPs. Unlike MLPs, which multiply the input by weights, KANs process the input through learnable B-spline functions [5]. In the intermediate layers, while MLPs use a fixed activation function as σ in $\sigma(w.x + b)$, KANs employ learnable activation functions based on parameterized B-spline functions, formulated as $\phi(x) = \text{silu}(x) + \sum_i c_i B_i(x)$, where B is the basis function and c is the trainable control point parameter. In this combination, the non-trainable component $\text{silu}(x)$ is added to the trainable spline component, transforming $\phi(x)$ into a residual activation function. Therefore, KAN consists of spline-based univariate functions along the network edges, designed as learnable activation functions. These features enhance both the accuracy and interpretability of the network. In our work, replacing the MLP head with KAN also enhances the accuracy of our model, as presented later in Section 4.2.

3 Proposed Method

In this paper, we aim to develop a lightweight model, focusing on both Multiply-Accumulate Operations (MACs) and the number of parameters. To achieve this,

Table 1: Summary of IQA datasets used in this work.

Database	# of Source Images	Dist. Type	# of Dist. Images	# of Dist. Types	Image Resolution
KONIQ-10K [15]	10,073	Authentic	10,073	–	1024×768px
SPAQ [9]	11,125	Authentic	11,125	–	Variable
BID [8]	–	Authentic	585	–	1280×960px to 2272×1704px
UHD-IQA [13]	6,073	Authentic	–	–	mostly 3840×2160px
KADID-10k [19]	81	Synthetic	10,125	25	512×384px
TID2013 [31]	25	Synthetic	3,000	24	512×384px
PIPAL [18]	250	Synthetic+Algorithmic	29,000	40	Variable

we utilize mobile architectures [22] in our image encoder. We propose a dual-branch architecture, comprising Authentic and Synthetic branches, as shown in Figure 1. Each branch includes a pre-processing module and a mobile-based image encoder. We use a KAN quality regression module to merge features from both branches and predict the final quality score. Initially, the two branches are treated as separate models and trained on authentic and synthetic datasets, respectively. Each branch employs a KAN regression head to independently predict image quality. After training the branches on their respective datasets, the image encoders from both the Authentic and Synthetic branches are integrated into the dual-branch model.

3.1 Datasets and Evaluation Metrics

In this work, we use seven publicly available IQA datasets: two with synthetic distortions (KADID-10K [19], TID2013 [31]), four with authentic distortions (KonIQ-10k [15], SPAQ [9], BID [8], and UHD-IQA [13]), the PIPAL [18] with both synthetic and algorithmic distortions. A summary of the datasets used in our experiments is provided in Table 1.

To evaluate the model’s performance, we use three common criteria: Spearman Rank Order Correlation Coefficient (SRCC) for prediction monotonicity, Pearson Linear Correlation Coefficient (PLCC) and Kendall Rank Correlation Coefficient (KRCC) for rank consistency.

3.2 Data Pre-processing

For each branch of our model, we applied distinct pre-processing methods during the training phase. In the Authentic branch, images were downsampled to 224x224 pixels for pretraining and 384x384 for UHD-IQA challenge, while in the Synthetic branch, images were kept at their original size, and only the cropping operation was applied. The rationale for downscaling in the Authentic branch is that the quality of Authentic datasets is intrinsically linked to their content (e.g., captured objects, image composition), and changes in image size do not impact perceived quality [36, 42]. However, in the Synthetic branch, downscaling can cause the loss of high-frequency details, which negatively impacts the predicted quality for images with synthetic distortions. Therefore, we only downscale images in the Authentic branch. For the Synthetic branch, we trained our model using cropped images of 224 sizes.

3.3 Quality Regression Module

For quality regression head we explore both MLPs and KANs. For KAN, we use the implementation provided by [1]. To ensure a fair comparison between the MLPs and KANs, we conducted two sets of comparisons. First, both models were evaluated using an identical architecture, consisting of two fully connected (FC) layers: the first layer reduced 1000 features to 128, followed by a second layer that reduced 128 features to 1. Given that KAN uses a higher number of parameters and FLOPs, we also adjusted the MLP model’s complexity to roughly match KAN’s FLOPs by adding an additional layer with 1125 neurons before the 128-neuron layer, using ReLU activation. The evaluations are conducted under two scenarios: within-dataset evaluation and cross-dataset evaluation, with further details discussed in Section 4.2.

3.4 Image Encoder Module

For image encoder module, we explore two lightweight CNN-based architectures and two ViT architectures. As demonstrated in [22], MobileNetV2 [34] and MobileNetV3 [16] are excellent choices for mobile CNN-based networks due to their reduced number of parameters and low computational complexity. Additionally, we consider the recent mobile ViT model, MobileCLIP-S2 [38], and MobileViT-S [25] known for its lightweight characteristics. It is important to note that although MobileCLIP-S2 is a VLM, in Section 4.3, we only utilize its vision encoder architecture, and our model does not follow a VLM training approach.

3.5 Loss Functions

MSE and PLCC (\mathcal{L}_{acc}). Distance-based losses, such as MAE and MSE, minimize the absolute difference between predicted scores and ground truth labels, ensuring accuracy in score prediction. In contrast, correlation-based losses like PLCC preserve the relative correlation of image quality, aligning better with human perception [6]. In this work, we enhance model accuracy by incorporating both distance-based (MSE) and correlation-based (PLCC) objectives as a combined loss function in all our experiments.

$$\mathcal{L}_{acc} = \alpha \cdot MSE + \beta \cdot PLCC \quad (1)$$

Color Space Robustness (\mathcal{L}_{rob}). As shown in traditional NSS-based models [4, 10, 23, 37], analyzing images in different color spaces improves the accuracy of visual quality assessments by providing deeper insights into perceptual attributes that affect image quality. The multi-color space approach captures diverse aspects of image quality that may not be covered by a single color space. Therefore, we propose a *color space loss* that evaluates image quality in the RGB, YUV, and LAB color spaces, resulting in more accurate and robust assessments. Our *color space loss* (\mathcal{L}_{rob}) is defined as

$$\widehat{MOS}_{\text{color-space}} = F_{\phi}(\mathbf{I}_{\text{color-space}}) \quad (2)$$

$$\sum_{\text{color-space}} \frac{1}{N} (\widehat{MOS}_{\text{color-space}} - MOS)^2 \quad (3)$$

where $\text{color-space} \in \{\text{RGB}, \text{YUV}, \text{LAB}\}$, and $\mathbf{I}_{\text{color-space}}$ represents the input image in the specified color space. Using our trained model F with parameters ϕ , we predict the MOS for $\mathbf{I}_{\text{color-space}}$, denoted by $\widehat{MOS}_{\text{color-space}}$. To compute the loss we calculate the mean squared error with the ground-truth MOS. By minimizing this *color space loss*, the model learns to align its output feature across various color-space representations. This approach enhances the model’s ability to predict image quality consistently, regardless of the color space used. It improves the model’s sensitivity to different distortions and ensures that assessments closely match human visual perception (see Section 4.4).

4 Evaluation Methodology and Results

In Section 4.1, we first discuss the implementation details of our model. We present the initial results of the Synthetic and Authentic models in Section 4.2 and explore replacing the MLP regression head with KAN. In Section 4.3, we experiment with different backbone networks as the image encoder for our Synthetic and Authentic models while keeping the regression head fixed as KAN. To enhance robustness, we investigate various loss functions in Section 4.4. Finally, we combine insights from these experiments and discuss the setup and integration of the Synthetic and Authentic branches in Section 4.5. In Section 4.6, we compare the performance of our final model against state-of-the-art methods based on accuracy metrics (SRCC, PLCC, KRCC, RMSE, and MAE) and computational complexity (MACs).

4.1 Implementation Details

The model training and testing is done on Pytorch on a single Nvidia A100 GPU. The models are trained for 100 epochs using the AdamW optimizer, starting with a learning rate of 5×10^{-5} and a weight decay of 1×10^{-4} . We use a custom scheduler, that includes a linear warmup phase followed by a cosine annealing schedule. The input image size varies depending on the branch: original size for the synthetic branch and resized 224x224 pixels for the authentic branch. For the UHD-IQA dataset, due to the large image resolutions, we resized images to 384x384 for the authentic branch and applied 1280x1280 center crops for the synthetic branch. The image encoder module in our final proposed model utilizes MobileNetV3, which is consistent across both branches. For the quality regression modules, we incorporate the KAN model to reduce the feature dimension and predict image quality, respectively.

4.2 Quality Regression Module: MLP versus KAN

In this experiment, we employed MobileNetV2 for training and utilized a standard loss function, which is a weighted combination of MSE and PLCC loss, measured between images in the batch. All other settings, such as the learning rate, remained constant. For the within-dataset evaluation, each dataset,

Table 2: Within-dataset evaluation results for different regression heads, MLP and KAN. The synthetic model is trained on KADID-10K and PIPAL datasets and evaluated using 10-fold cross-validation. Results indicate better performance of KAN compared to MLP on both datasets.

Dataset	KADID-10K			PIPAL		
	PLCC	SRCC	KRCC	PLCC	SRCC	KRCC
MLP	0.961	0.946	0.825	0.868	0.845	0.680
KAN	0.965	0.941	0.857	0.887	0.860	0.692

Table 3: Cross-dataset evaluation results for MLP and KAN as regression heads. Synthetic model is trained on the KADID-10K dataset and tested on the TID2013 dataset. Similarly, the Authentic model is trained on KONIQ-10K dataset and tested on BID dataset. Results indicate significantly better performance of KAN compared to MLP.

Model	Head			PLCC	SRCC	KRCC
	MLP	KAN	FLOPs			
KADID-10K / TID2013						
Synthetic	✓		0.26M	0.64	0.59	0.45
	✓		2.60M	0.67	0.63	0.48
		✓	2.31M	0.71	0.66	0.50
KONIQ-10K / BID						
Authentic	✓		0.26M	0.716	0.673	0.495
	✓		2.60M	0.726	0.673	0.501
		✓	2.31M	0.736	0.680	0.505

KADID-10K and PIPAL, were individually used for both training and testing using 10-fold cross-validation. The performance metrics are reported based on the average of the 10-fold cross-validation results. For the cross-dataset evaluation, we used two train/test pairs: KADID-10K / TID2013 for the Synthetic model and KONIQ-10K / BID for the Authentic model. The results of the within-dataset evaluation on KADID-10K and PIPAL are presented in Table 2. These results indicate that incorporating KAN as the regression head yields improvements over MLP. As previously discussed, we used two MLP regression heads: one with the same architecture and another with the same number of FLOPs. In the second experiment, we included the FLOPs for each MLP model and compared them to KAN. Table 3 presents the results of the cross-dataset evaluation, showing an improvement when using KAN compared to MLP, highlighting KAN’s superior generalizability. Hence, for the rest of the experiments, we only use KAN for quality regression task. It is important to note that making a fair comparison between KAN and MLP is challenging due to the differences in network topology and activation function choices, which significantly impact

Table 4: Comparing the performance of different image encoders. Among the evaluated encoders, MobileNetV3 consistently outperforms the others across all metrics.

Model	Backbone	#Parameters	PLCC	SRCC	KRCC
KADID-10K / TID2013					
Synthetic	MobileViT-S	5.6M	0.61	0.59	0.42
	MobileCLIP-S2 [38]	35.5M	0.69	0.65	0.46
	MobileNetV2 [34]	3.5M	0.71	0.66	0.50
	MobileNetV3 [16]	7.1M	0.72	0.68	0.50
KONIQ-10K / BID					
Authentic	MobileViT-S	5.6M	0.71	0.65	0.47
	MobileCLIP-S2	35.5M	0.78	0.75	0.57
	MobileNetV2	3.5M	0.74	0.68	0.51
	MobileNetV3	7.1M	0.79	0.78	0.60

MLP performance. Similarly, for KAN, factors such as the order of the spline and the number of spline intervals affect the results. Nevertheless, this paper demonstrates that incorporating KAN as part of the regression head can be as effective as using an MLP head.

4.3 Comparing Image Encoders

Using KAN as the regression head, we now compare different image encoders as the backbone network for our model. Similar to previous experiments, the Synthetic model is trained on KADID-10K and tested on TID2013, while the Authentic model is trained on KONIQ-10K and tested on BID. The performance metrics, including, PLCC, SRCC, and KRCC, are detailed in Table 4. The results show that both CNN-based models outperform the Mobile ViT models. The decreased performance of MobileCLIP-S2 can be attributed to the limited training data, as ViT-based models typically require large datasets to effectively learn the extensive network parameters. However, as reported in the next section, we observe that MobileCLIP-S2 helps in better generalization when using multiple color spaces. As expected, MobileNetV3 performs better than MobileNetV2; therefore, MobileNetV3 is used for the rest of our model design.

4.4 Comparing Various Loss Functions

Table 5 presents the performance of various loss functions considering different color spaces (RGB, YUV and LAB). The results indicate that incorporating both $MSE+PLCC$ loss (\mathcal{L}_{acc}) and the *color space loss* (\mathcal{L}_{rob}) improves the model’s performance. While the improvement in the RGB color space was not significant, major improvements were observed in the performance metrics for the YUV and LAB color spaces. The MobileCLIP-S2 [38] model, which uses the ViT architecture, is more robust and shows better performance across different color spaces. This could be attributed to the nature of ViT models, which normalize pixel values multiple times during processing.

Table 5: Performance comparison of various loss functions considering different color spaces. Best performing model is shown in Bold.

Model	Loss		Backbone	RGB			YUV			LAB		
	\mathcal{L}_{acc}	\mathcal{L}_{rob}		PLCC	SRCC	KRCC	PLCC	SRCC	KRCC	PLCC	SRCC	KRCC
KADID-10K / TID2013												
Synthetic	✓		Mobile-	0.72	0.68	0.50	0.41	0.34	0.25	0.47	0.39	0.28
	✓	✓	NetV3	0.73	0.70	0.53	0.49	0.47	0.35	0.55	0.52	0.37
Synthetic	✓		Mobile-	0.69	0.65	0.46	0.44	0.38	0.29	0.46	0.40	0.29
	✓	✓	CLIP-S2	0.69	0.66	0.48	0.58	0.54	0.39	0.59	0.55	0.40
KONIQ-10K / BID												
Authentic	✓		Mobile-	0.79	0.78	0.60	0.59	0.58	0.42	0.59	0.57	0.42
	✓	✓	NetV3	0.82	0.79	0.61	0.69	0.68	0.51	0.76	0.74	0.56
Authentic	✓		Mobile-	0.78	0.74	0.56	0.65	0.62	0.46	0.67	0.62	0.46
	✓	✓	CLIP-S2	0.78	0.77	0.58	0.74	0.72	0.53	0.74	0.73	0.54

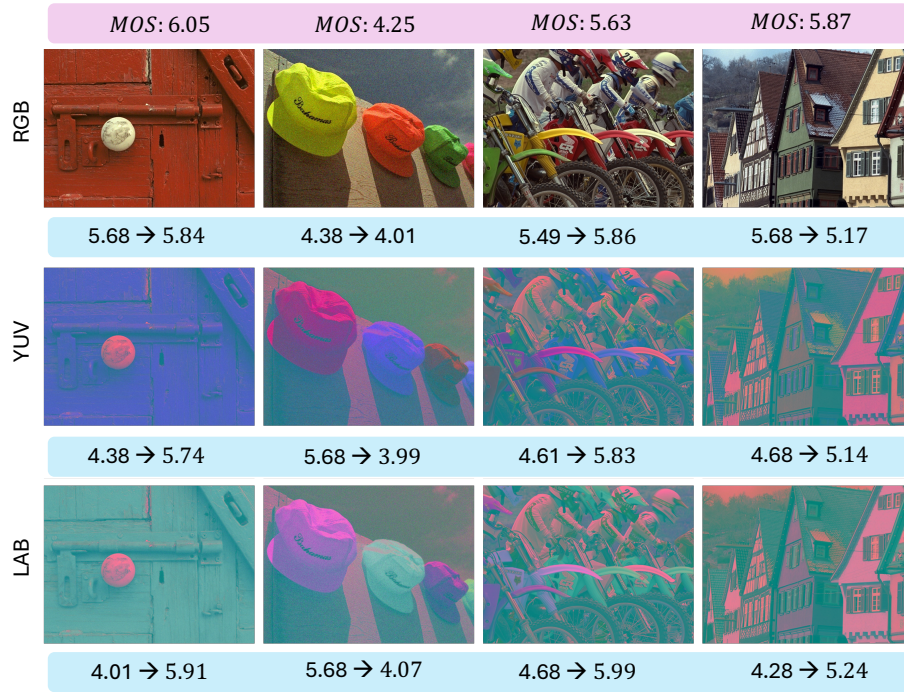


Fig. 2: MOS predictions of the Synthetic model using MobileNetV3, trained on KADID-10K and tested on sample images from the TID2013 dataset across different color spaces. Actual MOS scores are highlighted above in the pink colored bar. The change in predicted MOS scores before and after pre-training with *color space loss* (\mathcal{L}_{rob}) across various color spaces is shown in the blue colored bar. One can note that the predicted MOS scores for various images across different color spaces converge, indicating the robustness of the metric to different color spaces.

Table 6: Performance comparison of different training settings for LAR-IQA model on the validation set of the UHD-IQA dataset.

Branch		Finetune	KAN neurons	PLCC	SRCC	KRCC
Synthetic	Authentic					
✓			128	0.492	0.463	0.333
	✓		128	0.512	0.496	0.348
✓	✓		128	0.726	0.722	0.532
✓	✓		512	0.744	0.734	0.545
✓	✓	✓	512	0.809	0.803	0.611

Figure 2 shows the MOS predictions of the Synthetic model with MobileNetV3, trained on KADID-10K and tested on sample images from the TID2013 dataset. Despite the differences between the predicted MOS and ground truth values caused by varying scales and types of distortion in the two datasets, training with color space loss (\mathcal{L}_{rob}) enhances robustness and helps the model generalize better by ensuring that the results remain consistent across different color spaces.

4.5 Integrating Synthetic and Authentic Branches

In the previous section, we outlined the training process of the Synthetic branch using the KADID-10K dataset and the Authentic branch with the KONIQ-10K dataset. However, these datasets are limited in terms of number of images and distortions types. To enhance our model’s performance, we extended our training to include multiple datasets for each branch following a multi-task training process. Specifically, the Synthetic branch was trained on the KADID-10K, TID2013, and PIPAL datasets, while the Authentic branch was trained on the KONIQ-10K, BID, and SPAQ datasets. As outlined in Section 2.2, we mitigated subjective biases by integrating two branches trained on different distortion types, each with a distinct regression head for the respective datasets. We held out a random 10% of each dataset during the training process for the validation set and model selection per branch. We refer the reader to our previous paper [17], which details the multi-task training process for the IQA task.

The pretrained models are merged after removing the regression heads and adding new KAN heads. These new heads first downsample the embeddings of each branch, then concatenate the two branches, and finally add a regression head, as illustrated in Figure 1.

Training on UHD-IQA Dataset: We trained our two-branch model on UHD-IQA training dataset in three steps. First, we froze the two pre-trained branches (authentic and synthetic) and trained the KAN heads that reduced the output from 1000 to 256 (128 per branch) and then to a single output neuron. Second, we employed a larger KAN head with a 1000-dimensional embedding mapped to 512 per branch (1024 in total) and then to one neuron for the regression task. Finally, we fine-tuned the model using the pre-trained weights of the

Table 7: Evaluation of the performance of the baselines on the validation set of UHD-IQA. \uparrow means that higher values are better, \downarrow means that lower values are better. Best and second-best scores are highlighted in bold and underlined, respectively.

Method	PLCC \uparrow	SRCC \uparrow	KRCC \uparrow	RMSE \downarrow	MAE \downarrow	#Para \downarrow	MACs \downarrow
HyperIQA [35]	0.182	0.524	0.359	0.087	0.055	27.3M	<u>211G</u>
Effnet-2C-MLSP [41]	0.627	0.615	0.445	0.060	0.050	-	345G
CONTRIQUE [24]	0.712	0.716	0.521	<u>0.049</u>	<u>0.038</u>	27.9M	855G
ARNIQA [3]	0.717	0.718	0.523	0.050	0.039	27.9M	855G
CLIP-IQA+ [39]	0.732	0.743	0.546	0.108	0.087	102M	895G
QualiCLIP [2]	<u>0.752</u>	<u>0.757</u>	<u>0.557</u>	0.079	0.064	102M	901G
LAR-IQA (MLP head)	0.797	<u>0.791</u>	<u>0.601</u>	<u>0.042</u>	<u>0.033</u>	21.2M	$\leq 37G$
LAR-IQA (KAN head)	0.809	0.803	0.611	0.040	0.031	21.1M	$\leq 37G$

authentic and synthetic branches along with the larger KAN head on the new dataset.

Table 6 presents the results of various training strategies on the UHD-IQA dataset [13] used in the ECCV AIM challenge⁴. The model is trained using the training set and tested on the validation set. As observed, the model with the larger regression head, consisting of 512 layers, delivers improved performance. Furthermore, fine-tuning the model significantly improves the results, achieving a PLCC of approximately 0.81. Overall, the outcomes are promising. It is important to note that for the Authentic branch, we resized the images to 224x224x3, whereas for the Synthetic branch, we used the full-size image.

We present next the comparative evaluation of our proposed model compared to the baselines as evaluated on the validation and test set. The results for the baseline models are obtained from the original publication in [13].

4.6 Comparisons with the State-of-the-art

Table 7 and Table 8 presents a comparative evaluation result of our proposed model (LAR-IQA) compared to the SOTA models on the AIM UHD-IQA validation and test set respectively. It can be observed that our model LAR-IQA outperforms existing IQA models on both validation and test dataset in terms of all performance measures as well complexity measures. The proposed model is approx. $\times 5.7$ faster than the fastest SOTA model, HyperIQA. Furthermore, both branches together have around 21 million parameters, making it suitable for power and resource constrained mobile devices.

It should be noted that in both tables, the term "MLP head" refers to the MLP head with similar FLOPs to the KAN head. This aims to provide the reader with a performance comparison between KAN and MLP on the UHD-IQA dataset.

Additionally, we refer the reader to the UHD-IQA challenge [14], which compares various lightweight models proposed for UHD-IQA tasks. In the challenge

⁴ <https://codalab.lisn.upsaclay.fr/competitions/19335>

Table 8: Evaluation of the performance of the baselines on the test set of UHD-IQA. \uparrow means that higher values are better, \downarrow means that lower values are better. Best and second-best scores are highlighted in bold and underlined, respectively.

Method	PLCC \uparrow	SRCC \uparrow	KRCC \uparrow	RMSE \downarrow	MAE \downarrow	MACs(G) \downarrow
HyperIQA [35]	0.103	0.553	0.389	0.118	0.070	<u>211</u>
Effnet-2C-MLSP [41]	0.641	0.675	0.491	0.074	0.059	345
CONTRIQUE [24]	0.678	0.732	0.532	0.073	0.052	855
ARNIQA [3]	0.694	0.739	0.544	0.074	0.052	855
CLIP-IQA+ [39]	0.709	0.747	0.551	0.111	0.089	895
QualiCLIP [2]	0.725	0.770	0.570	0.083	0.066	901
LAR-IQA (MLP head)	<u>0.774</u>	<u>0.809</u>	<u>0.616</u>	0.058	<u>0.042</u>	\leq 37
LAR-IQA (KAN head)	0.787	0.836	0.642	<u>0.061</u>	0.041	\leq 37

rankings, LAR-IQA achieved second place. To ensure a fair comparison, we adhered to the challenge’s rules and dataset split.

5 Conclusion

To address the lack of low complexity, high accuracy and robust NR-IQA metric, we proposed in this work a lightweight NR-IQA model that surpasses the accuracy of current SOTA models while being more suitable for deployment on mobile devices. Our model proposes a dual-branch architecture that processes both authentic and synthetic distortions individually making it more robust to varied distortions as compared to models trained solely on single distortion types. Furthermore, we observe that training separately on multiple datasets helps to mitigate subjective biases inherent in individual datasets, thus improving further the model’s overall generalizability.

To ensure our model’s robustness across various color spaces, we incorporated RGB, YUV, and LAB into the training process. Furthermore, we explored KANs for the quality regression module instead of the commonly used MLPs. Our empirical results show that our proposed model not only surpasses SOTA NR-IQA models in accuracy but is also of much lower complexity, making it well-suited for real-world applications.

References

1. Efficient-kan implementation [source code], <https://github.com/Blealtan/efficient-kan.git>
2. Agnolucci, L., Galteri, L., Bertini, M.: Quality-aware image-text alignment for real-world image quality assessment. arXiv:2403.11176 (2024)
3. 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)
4. Bampis, C.G., Gupta, P., Soundararajan, R., Bovik, A.C.: Speed-qa: Spatial efficient entropic differencing for image and video quality. *IEEE signal processing letters* **24**(9), 1333–1337 (2017)
5. Bodner, A.D., Tepsich, A.S., Spolski, J.N., Pourteau, S.: Convolutional kolmogorov-arnold networks. arXiv:2406.13155 (2024)
6. Chen, Z., Wang, J., Li, B., Yuan, C., Hu, W., Liu, J., Li, P., Wang, Y., Zhang, Y., Zhang, C.: Gmc-iqa: Exploiting global-correlation and mean-opinion consistency for no-reference image quality assessment. arXiv:2401.10511 (2024)
7. Cheon, M., Yoon, S.J., Kang, B., Lee, J.: Perceptual image quality assessment with transformers. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 433–442 (2021)
8. Ciancio, A., da Silva, E.A., Said, A., Samadani, R., Obrador, P., et al.: No-reference blur assessment of digital pictures based on multifeature classifiers. *IEEE Transactions on image processing* **20**(1), 64–75 (2010)
9. Fang, Y., Zhu, H., Zeng, Y., Ma, K., Wang, Z.: Perceptual quality assessment of smartphone photography. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 3677–3686 (2020)
10. Ghadiyaram, D., Bovik, A.C.: Perceptual quality prediction on authentically distorted images using a bag of features approach. *Journal of vision* **17**(1), 32–32 (2017)
11. Golestaneh, S.A., Dadsetan, S., Kitani, K.M.: No-reference image quality assessment via transformers, relative ranking, and self-consistency. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 1220–1230 (2022)
12. He, C., Zheng, Q., Zhu, R., Zeng, X., Fan, Y., Tu, Z.: Cover: A comprehensive video quality evaluator. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops. pp. 5799–5809 (2024)
13. 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)
14. Hosu, V., Conde, M.V., Timofte, R., Agnolucci, L., Zadtootaghaj, S., Barman, N., et al.: AIM 2024 challenge on uhd blind photo quality assessment. In: Proceedings of the European Conference on Computer Vision (ECCV) Workshops (2024)
15. 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. 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)
17. Jamshidi Avanaki, N., Ghildiyal, A., Zadtootaghaj, S., Barman, N.: MSLIQA: Enhancing Learning Representations for Image Quality Assessment through Multi-Scale Learning. arXiv (2024)

18. Jinjin, G., Haoming, C., Haoyu, C., Xiaoxing, Y., Ren, J.S., Chao, D.: Pipal: a large-scale image quality assessment dataset for perceptual image restoration. In: *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI* 16. pp. 633–651. Springer (2020)
19. Lin, H., Hosu, V., Saupe, D.: Kadid-10k: A large-scale artificially distorted iqa database. In: *2019 Eleventh International Conference on Quality of Multimedia Experience (QoMEX)*. pp. 1–3. IEEE (2019)
20. Liu, L., Dong, H., Huang, H., Bovik, A.C.: No-reference image quality assessment in curvelet domain. *Signal Processing: Image Communication* **29**(4), 494–505 (2014)
21. Liu, Z., Wang, Y., Vaidya, S., Ruehle, F., Halverson, J., Soljačić, M., Hou, T.Y., Tegmark, M.: Kan: Kolmogorov-arnold networks (2024)
22. Luo, C., He, X., Zhan, J., Wang, L., Gao, W., Dai, J.: Comparison and benchmarking of ai models and frameworks on mobile devices. *arXiv:2005.05085* (2020)
23. Madhusudana, P.C., Birkbeck, N., Wang, Y., Adsumilli, B., Bovik, A.C.: St-greed: Space-time generalized entropic differences for frame rate dependent video quality prediction. *IEEE Transactions on Image Processing* **30**, 7446–7457 (2021)
24. 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)
25. Mehta, S., Rastegari, M.: Mobilevit: light-weight, general-purpose, and mobile-friendly vision transformer. *arXiv preprint arXiv:2110.02178* (2021)
26. Mittag, G., Zadtootaghaj, S., Michael, T., Naderi, B., Möller, S.: Bias-aware loss for training image and speech quality prediction models from multiple datasets. In: *2021 13th International Conference on Quality of Multimedia Experience (QoMEX)*. pp. 97–102. IEEE (2021)
27. Mittal, A., Moorthy, A.K., Bovik, A.C.: No-reference image quality assessment in the spatial domain. *IEEE Transactions on image processing* **21**(12), 4695–4708 (2012)
28. Mittal, A., Soundararajan, R., Bovik, A.C.: Making a “completely blind” image quality analyzer. *IEEE Signal processing letters* **20**(3), 209–212 (2012)
29. Moorthy, A.K., Bovik, A.C.: A two-step framework for constructing blind image quality indices. *IEEE Signal processing letters* **17**(5), 513–516 (2010)
30. Moorthy, A.K., Bovik, A.C.: Blind image quality assessment: From natural scene statistics to perceptual quality. *IEEE transactions on Image Processing* **20**(12), 3350–3364 (2011)
31. Ponomarenko, N., Jin, L., Ieremeiev, O., Lukin, V., Egiazarian, K., Astola, J., Vozel, B., Chehdi, K., Carli, M., Battisti, F., et al.: Image database tid2013: Peculiarities, results and perspectives. *Signal processing: Image communication* **30**, 57–77 (2015)
32. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: *International conference on machine learning*. pp. 8748–8763 (2021)
33. Saad, M.A., Bovik, A.C., Charrier, C.: Blind image quality assessment: A natural scene statistics approach in the dct domain. *IEEE transactions on Image Processing* **21**(8), 3339–3352 (2012)
34. Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., Chen, L.C.: Mobilenetv2: Inverted residuals and linear bottlenecks. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 4510–4520 (2018)

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: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3667–3676 (2020)
36. Sun, W., Min, X., Tu, D., Ma, S., Zhai, G.: Blind quality assessment for in-the-wild images via hierarchical feature fusion and iterative mixed database training. *IEEE Journal of Selected Topics in Signal Processing* **17**(6), 1178–1192 (2023)
37. Tu, Z., Yu, X., Wang, Y., Birkbeck, N., Adsumilli, B., Bovik, A.C.: Rapique: Rapid and accurate video quality prediction of user generated content. *IEEE Open Journal of Signal Processing* **2**, 425–440 (2021)
38. Vasu, P.K.A., Pouransari, H., Faghri, F., Vemulapalli, R., Tuzel, O.: Mobileclip: Fast image-text models through multi-modal reinforced training. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 15963–15974 (2024)
39. 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)
40. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image Quality Assessment: From Error Visibility to Structural Similarity. *IEEE Transactions on Image Processing* **13**(4), 600–612 (2004)
41. Wiedemann, O., Hosu, V., Su, S., Saupe, D.: Konx: cross-resolution image quality assessment. *Quality and User Experience* **8**(1), 8 (2023)
42. Wu, H., Liao, L., Chaofeng, C., Hou, J., Wang, A., Sun, W., Yan, Q., Lin, W.: Disentangling aesthetic and technical effects for video quality assessment of user generated content. [arXiv:2211.04894](https://arxiv.org/abs/2211.04894) (2022)
43. 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. In: International Conference on Machine Learning
44. Xu, J., Ye, P., Li, Q., Du, H., Liu, Y., Doermann, D.: Blind image quality assessment based on high order statistics aggregation. *IEEE Transactions on Image Processing* **25**(9), 4444–4457 (2016)
45. Xue, W., Mou, X., Zhang, L., Bovik, A.C., Feng, X.: Blind image quality assessment using joint statistics of gradient magnitude and laplacian features. *IEEE Transactions on Image Processing* **23**(11), 4850–4862 (2014)
46. Yang, S., Wu, T., Shi, S., Lao, S., Gong, Y., Cao, M., Wang, J., Yang, Y.: Maniqa: Multi-dimension attention network for no-reference image quality assessment. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1191–1200 (2022)
47. Ye, P., Kumar, J., Kang, L., Doermann, D.: Unsupervised feature learning framework for no-reference image quality assessment. In: 2012 IEEE conference on computer vision and pattern recognition. pp. 1098–1105. IEEE (2012)
48. Yun, Y.K., Lin, W.: You Only Train Once: A Unified Framework for Both Full-Reference and No-Reference Image Quality Assessment. [arXiv:2310.09560](https://arxiv.org/abs/2310.09560) (2024)
49. Zhang, L., Zhang, L., Bovik, A.C.: A feature-enriched completely blind image quality evaluator. *IEEE Transactions on Image Processing* **24**(8), 2579–2591 (2015)
50. Zhao, L., Wang, L.: A new lightweight network based on mobilenetv3. *KSIIT Transactions on Internet and Information Systems (TIIS)* **16**(1), 1–15 (2022)
51. 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:2405.19298](https://arxiv.org/abs/2405.19298) (2024)