# ATTENTION DOWN-SAMPLING TRANSFORMER, RELATIVE RANKING AND SELF-CONSISTENCY FOR BLIND IMAGE QUALITY ASSESSMENT

*Mohammed Alsaafin[a], Musab Alsheikh[b], Saeed Anwar[c,d], Muhammad Usman[c,d]*

King Fahd University of Petroleum and Minerals, Dhahran, 31261, Saudi Arabia
[a] Department of Industrial and Systems Engineering, [b] Department of Electrical Engineering, [c] Department of Information and Computer Science, [d] SDAIA-KFUPM Joint Research Center for Artificial Intelligence
{g201072600, g202114890, saeed.anwar, muhammad.usman}@kfupm.edu.sa

## ABSTRACT

The no-reference image quality assessment is a challenging domain that addresses estimating image quality without the original reference. We introduce an improved mechanism to extract local and non-local information from images via different transformer encoders and CNNs. The utilization of Transformer encoders aims to mitigate locality bias and generate a non-local representation by sequentially processing CNN features, which inherently capture local visual structures. Establishing a stronger connection between subjective and objective assessments is achieved through sorting within batches of images based on relative distance information. A self-consistency approach to self-supervision is presented, explicitly addressing the degradation of no-reference image quality assessment (NR-IQA) models under equivariant transformations. Our approach ensures model robustness by maintaining consistency between an image and its horizontally flipped equivalent. Through empirical evaluation of five popular image quality assessment datasets, the proposed model outperforms alternative algorithms in the context of no-reference image quality assessment datasets, especially on smaller datasets. Codes are available at https://github.com/mas94/ADTRS

***Index Terms***— No-Reference Image Quality Assessment, CNNs, Transformers, Self-Consistency, Relative Ranking
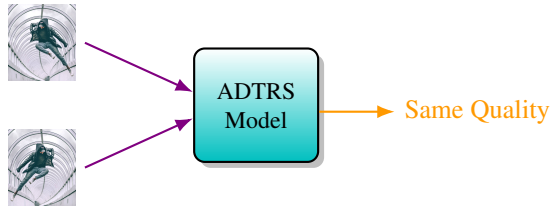
## 1. INTRODUCTION

Understanding image quality is essential for many applications; however, it may be difficult since we periodically need an ideal reference image. We can address this issue using Non-Reference Image Quality Assessment (NR-IQA). The objective is to develop techniques that can independently assess the image's quality without needing the original image. The importance of NR-IQA arises from its wide range of applications, including surveillance systems [1], medical imaging [2], content delivery networks [3], image & video compression [4], etc. It is vital in these domains to assess

quality without the original reference image. NR-IQA advances imaging technology and improves user experience. Existing NR-IQA methods focus on developing novel algorithms to handle the problem of evaluating image quality. Test Time Adaptation technique for Image Quality Assessment (TTAIQA) [5], Quality-aware Pre-Trained (QPT) [6] models through self-supervised learning, the Language-Image Quality Evaluator (LIQE), the data-efficient image quality transformer (DEIQT) [7] represents strides in this field and many methods that leverage CNNs. However, shortcomings persist, particularly the limitation imposed by the scarcity of labeled data, hindering the effectiveness of deep learning models and capturing only local features via CNNs while disregarding the nonlocal features of the image that transformers can capture. Popular datasets like the largest NR-IQA dataset, FLIVE, fall short compared to those in other domains, impeding the robust training of NR-IQA models.

Our main contribution is to develop an enhanced NR-IQA model to elevate its performance based on established metrics by leveraging the transformer architecture to capture nonlocal features and CNNs to capture local features. We seek to assess the performance of our improved model against existing NR-IQA methods. We test our methods using the most popular image quality datasets like LIVE, TID2013, CSIQ, LIVE-C and KonIQ10K. We'll utilize metrics like Spearman's Rank-Order Correlation Coefficient (SRCC) and Pearson's Linear Correlation Coefficient (PLCC) to calculate how well our model works. The main goal of this study is to get better-performing models using famous performance metrics.

## 2. RELATED WORKS

Image quality assessment, or IQA, is an integral part of computer vision and image processing and is widely utilized in social networking, online content sharing, and photography. This review is divided into two parts. The first discusses the recently released papers and focuses on developing creative methods for assessing NR-IRQA. The second part focuses

**Fig. 1**. Diagram illustrating the NR-IQA Model with inputs and outputs.

on papers that use convolutional neural networks (CNNs) and transformers for quality assessment.

Zhao et al. [6] provide a unique method for blind image assessment that uses self-supervised learning to get over the absence of labeled data. They propose a customized pretext task and a quality-aware contrastive loss, expanding image distortion simulations to enhance NR-IQA. Their Quality-conscious Pre-trained models beat current techniques on several BIQA benchmark datasets, demonstrating increased sensitivity for image quality. Distribution changes between training and testing situations in blind IQA significantly impact inference performance. Based on this, Subhadeep et al. [5] propose test-time adaptation by incorporating quality-relevant auxiliary activities at the batch and sample levels. Employing even a small batch of test data substantially improves model performance, surpassing existing SOTA approaches.

A multitask approach [7] for blind image quality evaluation, leveraging auxiliary information from other tasks and automating model parameter sharing and weighting for loss functions, benefiting from scene categorization and distortion type identification tasks. The mentioned method outperforms existing approaches across multiple IQA datasets, enhancing resilience and aligning better with quality annotations. It also offers insights into the creation of next-generation NR-IQA models.

Recently, a distinctive Mixture of Expert [8] introduces blind image quality assessment. This approach trains two separate encoders in an unsupervised setting to capture high-level content and low-level image characteristics. A linear regression model is developed to evaluate image quality by leveraging the synergy between these features. The authors showcase their technique's superiority over others across various extensive IQA databases, adeptly handling genuine and synthetic distortions. They highlight the significant impact of content-aware image representations, especially in an unsupervised context, on enhancing non-reference IQA performance.

The introduction of transformers [9] as novel network architecture defines a novel attention mechanism to draw global dependencies between I/O while eliminating the need for recurrence and convolutions. Furthermore, You et al. [10] were among the first papers to explore the transformer application in Image Quality (TRIQ) assessment. Building upon

the original Transformer encoder in Vision Transformer, they presented a shallow Transformer encoder atop a convolutional neural network (CNN)-extracted feature map. Their architecture accommodates images of varying resolutions, employing adaptive positional embedding. They systematically evaluated different Transformer configurations across multiple publicly available image quality databases. Their research outcomes highlight the efficacy of the proposed TRIQ framework.

For blind image quality assessment, Zhang et al. [11] presented a unique deep bilinear model that can handle real and artificial distortions. Two CNNs designed for various distortion conditions make up their model. One CNN is pre-trained on a sizable dataset for visual distortion classification to handle synthetic distortions, using a previously trained CNN to classify images for real distortions. The two CNNs' features were bilinearly pooled to provide a cohesive interpretation for image quality estimate. The performance of the entire model is improved by fine-tuning it on target subject-rated datasets using a variation of stochastic gradient descent. Numerous tests confirm the model's exceptional performance on artificial and real image databases to approach the issue of no reference IQA as a learning-to-rank problem in which training utilizes the ranking data.

CNNs are good at translating images but struggle with rotations. H-Nets [12], an innovative solution to the problem of imagining translation and rotation impacting computer vision tasks differently, can accomplish 360-degree rotation equivariance and patch-wise translation using circular harmonics instead of standard CNN filters. Each receptive field patch receives the maximum responsiveness and orientation from this special construction. Deep feature mappings inside the network may encode complicated rotational invariants thanks to H-Nets' parameter-efficient representation with constant computational complexity. Their approach is easily incorporated into contemporary systems such as batch normalization and deep supervision. H-Nets compete well on benchmark problems and produce cutting-edge classification scores on rotated-MNIST, demonstrating their effectiveness. Their research indicates that the data remains vulnerable to equivariant transformations even with various augmentation techniques to increase CNN generalization.
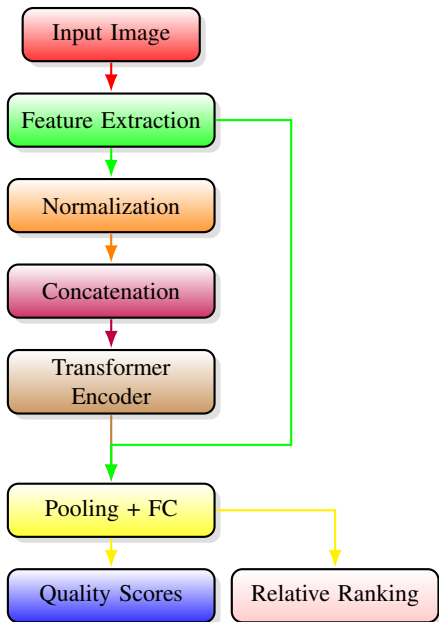
The field of blind image quality assessment has seen significant advancements recently, where each article addresses critical challenges in unique ways. The approaches leverage self-supervised learning, test-time adaptation, multitask learning, and unsupervised feature extraction to improve the precision and resilience of IQA models, making them better suited for real-world applications. Collectively, many studies contribute to the ongoing development of SOTA NR-IQA models, benefiting a wide range of industries and applications. Our primary goal in this paper centers around developing an advanced NR-IQA model that leverages transformers and CNNs to assess the quality of images with the pri-

mary aim of increasing its performance per established performance metrics. We intend to evaluate the ability of our enhanced model on five popular NR-IQA datasets by subjecting it to a thorough evaluation alongside existing NR-IQA approaches.

## 3. METHODOLOGY

In the proposed Attention Down-Sampling Transformer, Relative ranking and self-consistency abbreviated as ADTRS model for NR-IQA, we adopt the relative ranking and self-consistency mechanisms inspired by TReS [13] but employ a completely different transformer architecture to evaluate the quality of images without reference standards. Our method begins with an input image from which a series of CNN layers extract crucial features representing varying complexities and scales. These features are then normalized and subjected to dropout to ensure the model's generalizability across diverse image sets.

As depicted in Figure 2, the workflow progresses by concatenating extracted features to create a cohesive feature set. The Transformer encoder employs self-attention to emphasize essential data elements, leveraging attention mechanisms to prioritize relevant information. Subsequently, the encoder's output is aggregated through a fully connected layer, facilitating dimensionality reduction and regression analysis. By incorporating self-consistency mechanisms, the model ensures the reliability of its predictions. In the final stage of the ADTRS model, it produces both absolute quality scores and relative rankings, enabling comprehensive image quality assessment and facilitating meaningful comparisons.



**Fig. 2**. The basic building block of our the proposed ADTRS architecture.

### 3.1. Feature Extraction

For an input image $I$ defined in the space $\mathbb{R}^{3 \times m \times n}$ with dimensions $m$ and $n$ symbolizing the width and height, respectively, the objective is to evaluate its perceptual QS. A CNN, represented by $f_\phi$ with learnable parameters $\phi$, is utilized to extract features $F_i$ from the $i^{th}$ block, where $F_i \in \mathbb{R}^{b \times c_i \times m_i \times n_i}$ captures the feature maps with batch size $b$, and $c_i, m_i$, and $n_i$ denote the dimensions of the channels, width, and height of the extracted features, correspondingly.

In neural network architectures, specifically in the context of deep learning and CNNs, a series of pre-processing steps are often applied to the extracted features from different layers. These steps include normalization, pooling, and dropout, each serving distinct purposes to enhance the network's performance and generalization. Normalization is employed to address variations in feature scales among different layers. Standardizing the features to have zero mean and unit variance or scaling them to a specific range ensures that they contribute more uniformly to the learning process. This step is crucial because features with different scales might dominate the learning, hindering the network's ability to converge effectively.
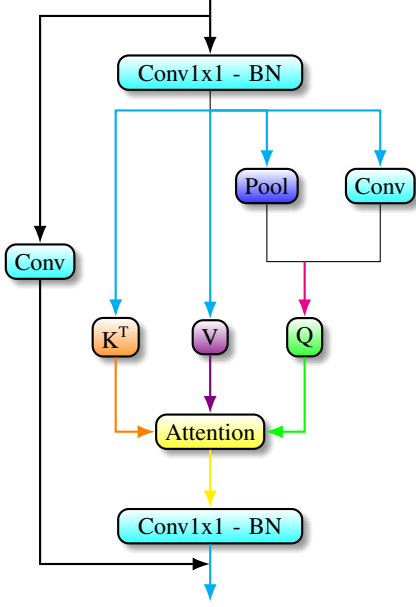
Pooling layers play a pivotal role in down-sampling the spatial dimensions of feature maps. By selecting the maximum or average values within specific regions, pooling reduces the computational complexity of the network, making it more efficient. Additionally, pooling contributes to the network's robustness by detecting invariant features and handling spatial variations in the input data. As a regularization technique, dropout addresses the risk of over-fitting by randomly deactivating a fraction of neurons during training. By preventing the network from relying too heavily on specific neurons, dropout encourages learning more robust features and improves the model's ability to generalize to unseen data. During testing, all neurons are reinstated to ensure the full utilization of the trained network. These pre-processing steps: normalization, pooling, and dropout—collectively constitute a comprehensive strategy for enhancing the efficiency, generalization, and robustness of neural networks, particularly in the complex tasks associated with deep learning architectures. (Eq. 1) is used to normalize the feature vector $F_i$ using the Euclidean norm. The $L_2$ pooling is defined by (Eq. 2).

$$F_i = \frac{F_i}{\max \left( \|F_i\|_2, \epsilon \right)} \tag{1}$$

$$P(x) = \sqrt{g * (x \odot x)} \tag{2}$$

where $\odot$ denotes the point-wise product, and the blurring kernel $g(\cdot)$ is implemented via a Hamming window that approximately applies the Nyquist criterion.

The extracted features will be concatenated after going through the normalization, pooling and dropout layers.
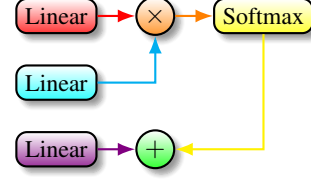
**Fig. 3**. Dual-Path downsampling transformer encoder [15] adopted in our proposed method.

## 3.2. Transformer Encoder

We adopt the encoder architecture from the literature to process multi-scale features $F_i$ from the CNN [14]. These features are sequenced and fed into the Transformer encoder [15], where a multi-head, multi-layer self-attention mechanism, depicted in Figure 2, is employed to model the dependencies across the feature maps. The Transformer's architecture can learn complicated feature relationships without any built-in inductive bias.

At the heart of our model lies the Transformer Encoder Layer paired with the Self-attention mechanism, essential for analyzing image characteristics through sophisticated spatial recognition. Initially, the model processes image features that distill the essence of the visual input and integrates positional encoding to inject spatial context into these features. Following this, the multi-head self-attention framework comes into play, dissecting and assessing input segments by generating attention scores through queries, keys, and values, all created via adaptive linear transformations.

After the self-attention phase, the outputs are fused and normalized in the Add & Norm step, a measure that stabilizes the learning process and embeds residual connections, vital for the architecture's depth and efficacy. A subsequent Feed-Forward Network (FFN) executes additional linear transformations, punctuated by ReLU activation, to polish the feature set further. The depth of this encoding process is represented by 'Nx', reflecting the iterations of the transformation sequence. The depicted data flow in Figure 3, especially the loopback arrows, emphasizes the Transformer's residual connections, ensuring a seamless and continuous information



**Fig. 4**. Schematic of the Self-Attention Mechanism.

stream within the model's architecture.

The self-attention mechanism within the Transformer model, as depicted in Figure 4, commences with dual linear blocks that reformulate the input data into intermediate states, integral for the ensuing attention calculations. This mechanism further employs the SoftMax function to normalize attention scores calculated from the dot products of queries and keys. This normalization facilitates a targeted distribution of attention across different data segments. Following this, the output from the SoftMax stage is combined with the outputs of a third linear block, symbolizing information integration within the attention process. This synthesis is not terminal but instead feeds back iteratively into the SoftMax stage, highlighting the recursive nature of the self-attention mechanism. This iterative loop, crucial for refining attention over successive cycles, is visually represented in Figure 4.

**Multi-Head Attention**: The multi-head attention mechanism, pivotal in our model, involves transforming input features into query, key, and value vectors. These vectors are then processed through the attention mechanism as shown in the following equations:

$$\text{MultiHead}\left(Q', K', V'\right) = \text{Concat}\left(\mathrm{h}_1, \ldots, \mathrm{h}_h\right) W^o, \quad (3)$$

$$\mathrm{h}_i = \text{Attention}\left(Q_i, K_i, V_i\right), \quad (4)$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) \odot V, \quad (5)$$

where $W^o$ is a projection matrix, $Q$, $K$, and $V$ refer to the Query, Key, and Value matrices, respectively, and the softmax attention is normalized over the key dimension $d_k$.

## 3.3. Feature Fusion and Quality Estimation

We utilize fully connected layers to combine the extracted features from convolutional and self-attention mechanisms. This fusion is instrumental in predicting the image's quality (as visualized in Figure 2). Formulated as follows, the model is trained to minimize the regression loss as:

$$\mathcal{L}_{\mathrm{Q},B} = \frac{1}{N} \sum_{i=1}^{N} \|q_i - s_i\|, \quad (6)$$

where $q_i$ represents the predicted quality score for the $i^{th}$ image, while $s_i$ represents ground truth quality score.

**Table 1**. Summary of IQA datasets where "Dist." stands for Distortion, and "#" is for number.

| Databases | # of Dist. Images | # of Dist. Types | Dist. Type |
|---|---|---|---|
| LIVE | 799 | 5 | Synthetic |
| CSIQ | 866 | 6 | Synthetic |
| LIVE-C | 1.162 | - | Real |
| TID2013 | 3,000 | 24 | Synthetic |
| KonIQ10K | 10,073 | - | Real |

### 3.4. Relative Ranking Incorporation

The regression loss effectively handles quality prediction; it overlooks ranking and correlation among images. We aim to account for the relative ranking within batches, focusing on extreme cases due to computational constraints. In image batch $B$, $qa_{max}$, $qa'_{max}$, $qa_{min}$, and $qa'_{min}$ represent predicted qualities for the highest, second highest, lowest, and second lowest subjective quality scores, respectively. Utilizing triplet loss with $d(x, y) = |x - y|$, we aim for constraints like $d(qa_{max}, qa'_{max}) + \text{margin}_1 \leq d(qa_{max}, qa_{min})$. Similarly, we desire $d(qa_{min}, qa'_{min}) + \text{margin}_2 \leq d(qa_{max}, qa_{min})$. Empirically selecting margin values is challenging due to dataset variations. For perfect predictions, $\text{margin}_1$ is bounded by $sqa'_{max} - sqa_{min}$, serving as an upper-bound during training, where $\text{margin}_1 = sqa'_{max} - sqa_{min}$, where $sqa'_{max}$ signifies the subjective quality score associated with the image having the predicted quality score $qa'_{max}$. We can do a similar process for $\text{margin}_2$ which is bounded by $sqa_{max} - sqa'_{min}$.

$$
\begin{aligned}
\mathcal{L}_{\text{RR},B} = & \\
\mathcal{L}_{\text{triplet}} &\left( qa_{max}, qa'_{max}, qa_{min} \right) + \mathcal{L}_{\text{triplet}} \left( qa_{min}, qa'_{min}, qa_{max} \right) \\
= & \max \left\{ 0, d\left( qa_{max}, qa'_{max} \right) - d\left( qa_{max}, qa_{min} \right) + \text{margin}_1 \right\} \\
& + \max \left\{ 0, d\left( qa'_{min}, qa_{min} \right) - d\left( qa_{max}, qa_{min} \right) + \text{margin}_2 \right\}
\end{aligned}
\tag{7}
$$

### 3.5. Self-Consistency Mechanism

For the last part of the methodology, we advocate leveraging the model's uncertainty in both the original input image and its equivariant transformation during the training process. To enhance the robustness of the model, we exploit self-consistency by establishing a self-supervisory signal between each image and its equivariant transformation. For a given input $I$, denote the output logits from the Convolutional and Transformer layers as $\zeta_{\epsilon,\text{conv}}(I)$ and $\zeta_{\psi,\text{atten}}(I)$, respectively, where $\zeta_{\epsilon,\text{conv}}$ and $\zeta_{\psi,\text{atten}}$ represent the CNN and Transformer with learnable parameters $\epsilon$ and $\psi$, respectively. Our model utilizes these outputs to predict image quality. Given that human subjective scores remain consistent for the horizontally flipped version of the input image, we anticipate $\zeta_{\epsilon,\text{conv}}(I) = \zeta_{\epsilon,\text{conv}}(\tau(I))$ and $\zeta_{\psi,\text{atten}}(I) = \zeta_{\psi,\text{atten}}(\tau(I))$, where $\tau$ signifies the horizontal flipping transformation. Consequently, by incorporating our consistency loss, the network learns to fortify its representation learning autonomously,

eliminating the need for additional labels or external supervision. We aim to minimize the self-consistency loss

$$
\begin{aligned}
\mathcal{L}_{\text{SC}} = & \|\zeta_{\epsilon,\text{conv}}(I) - \zeta_{\epsilon,\text{conv}}(\tau(I))\| + \\
& \|\zeta_{\psi,\text{atten}}(I) - \zeta_{\psi,\text{atten}}(\tau(I))\| + \theta_1 \left\| \mathcal{L}_{\text{RR},B} - \mathcal{L}_{\text{RR},\tau(B)} \right\|,
\end{aligned}
\tag{8}
$$

where $\tau(B)$ signifies the equivariant transformation on image batch $B$.

### 3.6. Composite Loss Function

The overall training process involves the minimization of a composite loss function (Eq. 9), which encompasses quality loss (Eq. 6), relative ranking loss (Eq. 7), and self-consistency loss (Eq. 8). These losses are balanced by coefficients $\theta_1$, $\theta_2$, and $\theta_3$ to optimize the training outcome effectively. Our proposed model aims to provide a comprehensive and reliable NR-IQA solution by incorporating these elements. The effectiveness of this approach is demonstrated in the experimental results section.

$$
\mathcal{L}_{\text{CLF}} = \mathcal{L}_{\text{Q}} + \theta_2 \mathcal{L}_{\text{RR}} + \theta_3 \mathcal{L}_{\text{SC}}
\tag{9}
$$

## 4. EXPERIMENTS

We assess our proposed ADTRS model's performance on five widely recognized IQA datasets, displayed in Table 1 (among the distortions, three were synthetically generated while two occurred authentically). We utilize two standard performance metrics, PLCC and SROCC, among various metrics available in IQA evaluation. PLCC (Pearson Linear Correlation Coefficient) evaluates the correlation between algorithmic results and human eye subjective scores, reflecting the algorithm's accuracy. On the other hand, SROCC (Spearman Rank-Ordered Correlation Coefficient) measures the monotonicity of the algorithm's predictions. Both metrics range from 0 to 1, with higher values indicating superior performance.

**Implementation Details**: We used an NVIDIA RTX 2060 GPU and PyTorch to train our model for training and testing. We augmented the horizontal and vertical dimensions of 138 randomly chosen patches, each measuring 224 by 224 pixels, from each image, by accepted IQA training protocols. The quality scores of the original image were carried over to these patches. With a weight decay of $5 \times 10^{-4}$ across a maximum of five epochs, we trained by minimizing the composite loss function over the training set, changing the learning rate from $2 \times 10^{-5}$ and decreasing it by a factor of 5 after each epoch. A total of 138 patches, each measuring $224 \times 224$, were chosen randomly from the test picture during testing, and the final quality score was calculated by averaging their projected values. ResNet50 [28] served as the CNN backbone for our model, which was seeded using Imagenet weights. We set the hyperparameters $\theta_1$, $\theta_2$, and $\theta_3$ to $0.5, 0.05$, and $1$ accordingly, using the Transformer architecture with 4 encoder layers, a hidden layer dimensionality of 64 ($d = 16$), and

**Table 2**. Comparison of ADTRS and No-Reference State-of-the-Art Algorithms.

| | LIVE | | CSIQ | | TID2013 | | LIVE-C | | KonIQ10K | |
|---|---|---|---|---|---|---|---|---|---|---|
| | PLCC | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC | SROCC |
| DBCNN [11] | 0.971 | 0.968 | 0.959 | **0.946** | 0.865 | 0.816 | 0.869 | **0.869** | 0.884 | 0.875 |
| TReS [13] | 0.968 | 0.969 | 0.942 | 0.922 | 0.883 | 0.863 | 0.877 | 0.846 | **0.928** | **0.915** |
| HFD [16] | 0.971 | 0.951 | 0.890 | 0.842 | 0.681 | 0.764 | - | - | - | - |
| PQR [17] | 0.971 | 0.965 | 0.901 | 0.873 | 0.864 | 0.849 | 0.836 | 0.808 | - | - |
| DIIV INE [18] | 0.908 | 0.892 | 0.776 | 0.804 | 0.567 | 0.643 | 0.591 | 0.588 | 0.558 | 0.546 |
| BRISQUE [19] | 0.944 | 0.929 | 0.748 | 0.812 | 0.571 | 0.626 | 0.629 | 0.629 | 0.685 | 0.681 |
| ILNIQE [20] | 0.906 | 0.902 | 0.865 | 0.822 | 0.648 | 0.521 | 0.508 | 0.508 | 0.537 | 0.523 |
| BIECON [21] | 0.961 | 0.958 | 0.823 | 0.815 | 0.762 | 0.717 | 0.613 | 0.613 | 0.654 | 0.651 |
| MEON [22] | 0.955 | 0.951 | 0.864 | 0.852 | 0.824 | 0.808 | 0.710 | 0.697 | 0.628 | 0.611 |
| WaDIQaM [23] | 0.955 | 0.960 | 0.844 | 0.852 | 0.855 | 0.835 | 0.671 | 0.682 | 0.807 | 0.804 |
| TIQA [24] | 0.965 | 0.949 | 0.838 | 0.825 | 0.858 | 0.846 | 0.861 | 0.845 | 0.903 | 0.892 |
| MetaIQA [25] | 0.959 | 0.960 | 0.908 | 0.899 | 0.868 | 0.856 | 0.802 | 0.835 | 0.856 | 0.887 |
| P2P-BM [26] | 0.958 | 0.959 | 0.902 | 0.899 | 0.856 | 0.862 | 0.842 | 0.844 | 0.885 | 0.872 |
| HyperIQA [27] | 0.966 | 0.962 | 0.942 | 0.923 | 0.858 | 0.840 | **0.882** | 0.859 | 0.917 | 0.906 |
| ADTRS (Ours) | **0.972** | **0.970** | **0.960** | 0.943 | **0.897** | **0.878** | 0.864 | 0.836 | 0.918 | 0.905 |

16 heads ($h = 16$). All experiments were carried out with a consistent setup, adhering to NR-IQA standards. Datasets were randomly divided into $80\%$ and $/20\%$ train/test ratios.

**Performance Evaluation**: Table 2 presents a comprehensive performance comparison based on PLCC and SROCC metrics across five standard image quality datasets. Our ADTRS consistently exhibits superior performance in both PLCC and SROCC evaluations compared to existing algorithms. Notably, among these SOTA NR-IQA algorithms, some incorporated CNNs while others did not. We chose the following algorithms since they were the most recent algorithms in [13] based on the same performance metrics. The bolded entries in Table 2 are the best-performing algorithms for each dataset. Drawing inspiration from TReS [13], we adopted the relative ranking and self-consistency mechanism, making it equitable to compare our improved model (ADTRS) with it. As evident from Table 2, our model, employing a distinct transformer architecture, outperforms TReS and all other algorithms, particularly on smaller/synthetic datasets. Our model excels in Live, TID2013, and CSIQ's PLCC when compared against all different algorithms. Furthermore, our model performs exceptionally on real distorted and larger datasets, such as LIVE-C and KonIQ10K. In the KonIQ10K dataset, our model stands as the second-best performer, with TReS leading.

**Ablation Study & Hyperparameter Tuning**: As previously mentioned, we employed ResNet50 as the experiments' primary backbone; smaller backbones offered faster processing speeds, yielding relatively inferior results. We also explored ResNet34 and ResNet18, but they produced significantly worse outcomes. Moreover, we experimented with varying sample patch sizes for training and testing hyperparameters, ranging from 16 to 138. The results represent the optimal outcomes from our extensive experiments, achieved with 138 sample patches. Additionally, we investigated en-

hancing the batch size from 8 to 16, which led to poorer results than using 8 batches. To fine-tune our model, we conducted a grid search across hyperparameters. This involved running the grid search for 5 epochs and adjusting parameters with significant performance impacts. For instance, we explored different configurations of encoder layers, such as 2, 4, 6, and 8 layers, along with varying values for other influential hyperparameters, e.g., training and testing sample patches.

## 5. CONCLUSION

Our study presents an enhanced NR-IQA algorithm that efficiently merges CNNs and Transformer features, exploiting both local and non-local image characteristics for a comprehensive representation. We have incorporated a relative ranking loss function to capture essential ranking information among images, thereby augmenting the discriminative power of our model. By utilizing equivariant image transformations for self-supervision, we have bolstered the robustness of our approach. The performance of our method across five distinct IQA datasets underscores its robustness and adaptability. Our proposed algorithm outperforms all other algorithms on smaller and synthetic datasets while performing exceptionally well on larger datasets compared to different state-of-the-art algorithms. The results unequivocally establish the robustness and precision of our proposed method compared to the TReS model in accurately assessing image quality, highlighting its significant potential for diverse applications in image analysis and assessment.

# 6. REFERENCES

[1] W. Lu, W. Sun, X. Min, Z. Zhang, T. Wang, W. Zhu, X. Yang, and G. Zhai, "Blind surveillance image quality assessment via deep neural network combined with the visual saliency," in *ICAI*, Springer, 2022. 1

[2] S. Li, J. He, Y. Wang, Y. Liao, D. Zeng, Z. Bian, and J. Ma, "Blind ct image quality assessment via deep learning strategy: initial study," in *Image Perception, Observer Performance, and Technology Assessment*, SPIE, 2018. 1

[3] G. Yue, C. Hou, W. Yan, L. K. Choi, T. Zhou, and Y. Hou, "Blind quality assessment for screen content images via convolutional neural network," *DSP*, 2019. 1

[4] R. Hu, Y. Liu, Z. Wang, and X. Li, "Blind quality assessment of night-time image," *Displays*, 2021. 1

[5] S. Roy, S. Mitra, S. Biswas, and R. Soundararajan, "Test time adaptation for blind image quality assessment," in *ICCV*, 2023. 1, 2

[6] K. Zhao, K. Yuan, M. Sun, M. Li, and X. Wen, "Quality-aware pre-trained models for blind image quality assessment," in *CVPR*, 2023. 1, 2

[7] W. Zhang, G. Zhai, Y. Wei, X. Yang, and K. Ma, "Blind image quality assessment via vision-language correspondence: A multitask learning perspective," in *CVPR*, 2023. 1, 2

[8] A. Saha, S. Mishra, and A. C. Bovik, "Re-iqa: Unsupervised learning for image quality assessment in the wild," in *CVPR*, 2023. 2

[9] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *NIPS*, 2017. 2

[10] J. You and J. Korhonen, "Transformer for image quality assessment," in *ICIP*, IEEE, 2021. 2

[11] W. Zhang, K. Ma, J. Yan, D. Deng, and Z. Wang, "Blind image quality assessment using a deep bilinear convolutional neural network," *TCSVT*, 2018. 2, 6

[12] D. E. Worrall, S. J. Garbin, D. Turmukhambetov, and G. J. Brostow, "Harmonic networks: Deep translation and rotation equivariance," in *CVPR*, 2017. 2

[13] S. A. Golestaneh, S. Dadsetan, and K. M. Kitani, "No-reference image quality assessment via transformers, relative ranking, and self-consistency," in *WACV*, 2022. 3, 6

[14] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *ECCV*, Springer, 2020. 4

[15] Y. Li, J. Hu, Y. Wen, G. Evangelidis, K. Salahi, Y. Wang, S. Tulyakov, and J. Ren, "Rethinking vision transformers for mobilenet size and speed," in *ICCV*, 2023. 4

[16] J. Wu, J. Zeng, Y. Liu, G. Shi, and W. Lin, "Hierarchical feature degradation based blind image quality assessment," in *ICCV Workshops*, 2017. 6

[17] H. Zeng, L. Zhang, and A. C. Bovik, "A probabilistic quality representation approach to deep blind image quality prediction," *arXiv*, 2017. 6

[18] M. A. Saad, A. C. Bovik, and C. Charrier, "Blind image quality assessment: A natural scene statistics approach in the dct domain," *TIP*, 2012. 6

[19] A. Mittal, A. K. Moorthy, and A. C. Bovik, "No-reference image quality assessment in the spatial domain," *TIP*, 2012. 6

[20] L. Zhang, L. Zhang, and A. C. Bovik, "A feature-enriched completely blind image quality evaluator," *TIP*, 2015. 6

[21] J. Kim and S. Lee, "Fully deep blind image quality predictor," *JSTSP*, 2016. 6

[22] K. Ma, W. Liu, K. Zhang, Z. Duanmu, Z. Wang, and W. Zuo, "End-to-end blind image quality assessment using deep neural networks," *TIP*, 2017. 6

[23] S. Bosse, D. Maniry, K.-R. Müller, T. Wiegand, and W. Samek, "Deep neural networks for no-reference and full-reference image quality assessment," *TIP*, 2017. 6

[24] F. Gao, D. Tao, X. Gao, and X. Li, "Learning to rank for blind image quality assessment," *TNNLS*, 2015. 6

[25] H. Zhu, L. Li, J. Wu, W. Dong, and G. Shi, "Metaiqa: Deep meta-learning for no-reference image quality assessment," in *CVPR*, 2020. 6

[26] Z. Ying, H. Niu, P. Gupta, D. Mahajan, D. Ghadiyaram, and A. Bovik, "From patches to pictures (paq-2-piq): Mapping the perceptual space of picture quality," in *CVPR*, 2020. 6

[27] S. Su, Q. Yan, Y. Zhu, C. Zhang, X. Ge, J. Sun, and Y. Zhang, "Blindly assess image quality in the wild guided by a self-adaptive hyper network," in *CVPR*, 2020. 6

[28] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016. 5