# CLIP-Guided Attribute Aware Pretraining for Generalizable Image Quality Assessment

Daekyu Kwon[1], Dongyoung Kim[1], Sehwan Ki[2], Younghyun Jo[2], Hyong-Euk Lee[2], and Seon Joo Kim[1]

[1] Yonsei University
[2] Samsung Advanced Institue of Technology

**Abstract.** In no-reference image quality assessment (NR-IQA), the challenge of limited dataset sizes hampers the development of robust and generalizable models. Conventional methods address this issue by utilizing large datasets to extract rich representations for IQA. Also, some approaches propose vision language models (VLM) based IQA, but the domain gap between generic VLM and IQA constrains their scalability. In this work, we propose a novel pretraining framework that constructs a generalizable representation for IQA by selectively extracting quality-related knowledge from VLM and leveraging the scalability of large datasets. Specifically, we carefully select optimal text prompts for five representative image quality attributes and use VLM to generate pseudo-labels. Numerous attribute-aware pseudo-labels can be generated with large image datasets, allowing our IQA model to learn rich representations about image quality. Our approach achieves state-of-the-art performance on multiple IQA datasets and exhibits remarkable generalization capabilities. Leveraging these strengths, we propose several applications, such as evaluating image generation models and training image enhancement models, demonstrating our model's real-world applicability. We will make the code available for access.

**Keywords:** Image Quality Assessment · Visual Language Model · Real-world Application

## 1 Introduction

No-reference image quality assessment (NR-IQA) [10,25,40,41,45,53] is a task of quantifying the quality of images without a pristine reference image. Recently, methods in IQA have also started incorporating deep learning [2,15,21,47,49,50], similar to other fields in computer vision. However, effective application of deep learning in image quality assessment (IQA) faces challenges due to the limited size of existing IQA datasets [6,9,13,30,46]. Training an IQA model from scratch with a small dataset encounters difficulties in learning rich representations for image quality. This often results in degraded performance and poor generalization, thereby restricting the practical applicability of IQA in real-world scenarios.

To address the generalization issue derived from limited dataset size, IQA approaches have been developed to leverage rich representations from large
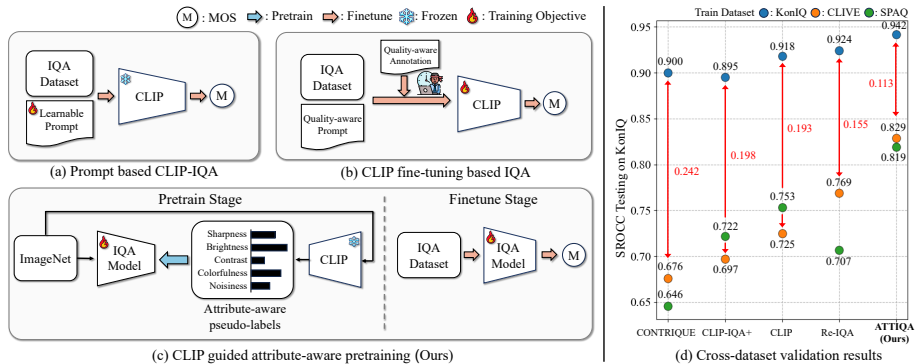
**Fig. 1:** An illustration of a training strategy for IQA across previous works and ours. (a) CLIP-IQA directly utilizes CLIP to predict MOS using zero-shot inference or prompt optimization. (b) LIQE fine-tunes CLIP for IQA using additional quality annotation, which requires human labor. (c) Our method incorporates the rich representation of a large dataset and CLIP's capability for IQA. We pretrain IQA model with attribute-aware pseudo-labels derived from CLIP and fine-tune it to the target IQA dataset. (d) Cross-dataset validation results, obtained by testing on the KonIQ dataset after training on various datasets. ATTIQA not only achieves state-of-the-art result, but also demonstrates superior generalization capability on unseen datasets, exhibiting less performance decline on cross-dataset setup compared to other methods.

datasets [34]. In [1, 41, 44, 50], transfer learning strategy was utilized by pretraining a model on ImageNet [34]. Several studies [21, 23, 36, 50, 52] have proposed IQA-specific pretext tasks, founded on the premise that distortions in images directly impact their quality. In this line of research, they aim to leverage the scalability of large datasets, developing methods that effectively extract quality-aware representation in a self-supervised learning manner.

Recently, Vision Language Models (VLM), exemplified by CLIP [31] have emerged as a robust framework in computer vision, highlighting their adaptability and generalization capabilities. Building on these strengths, exploiting VLM for IQA has also been explored. CLIP-IQA [43] introduced zero-shot IQA using the quality-aware prompt, demonstrating the applicability of CLIP for IQA (Fig. 1(a)). However, as CLIP is trained under the generic image-text pairs, there remains a need for research on how to apply CLIP more effectively within the IQA domain. In this regard, LIQE [51] directly adapts CLIP to IQA using text prompts about distortion and scene information(Fig. 1(b)). While LIQE effectively enhances CLIP's representation for IQA, this strategy is constrained by the necessity of supplementary annotations (scene and distortion type) for direct CLIP training, which requires additional human labor.

In this work, we introduce a method that exhibits enhanced generalization capabilities by effectively incorporating CLIP's extensive knowledge and the scalability of a large dataset. Our approach begins with the observation of previous works [43, 51] that CLIP inherently contains representations relevant to IQA. Building on this insight, we emphasize utilizing its inherent knowledge and pro-

pose an effective strategy to specialize CLIP for IQA by selectively extracting quality relevant information from its representation. To implement this strategy, we develop semi-supervised pretraining methods guided by the zero-shot inference of CLIP. Rather than training CLIP itself, we employ it to generate pseudo-labels from unlabeled datasets with quality-aware prompts, which are then used to train a target encoder(Fig. 1(c)). Moreover, as our approach does not require extra annotations, we can incorporate this approach with an arbitrary large dataset. Therefore, we can effectively transfer CLIP's quality-related knowledge with the scalability benefits of large datasets into the target encoder, resulting in the construction of a robust quality-aware representation.

Additionally, instead of relying on generic prompts such as "a good photo / a bad photo" as used in CLIP-IQA [43], our method incorporates more specific prompts related to distinct image attributes. These prompts are carefully selected through our specially designed prompt selection strategy. With these tailored prompts, we employ CLIP to generate pseudo-labels for a wide range of image attributes on a large-scale dataset. Recognizing the inherent challenges in quantifying image attributes, we adopt a ranking-based loss for training our IQA model. Taking advantage of using a large-scale dataset combined with the novel image attribute-aware pretraining strategy with CLIP guidance, our IQA framework significantly enhances the learning of representations closely associated with image attributes and quality. We propose "**ATTIQA**", **ATT**ribute-aware **IQA** with CLIP guided pretraining, achieving state-of-the-art performance on multiple IQA in the wild datasets, as well as on an aesthetic quality dataset.

For IQA, the ability to generalize beyond the training dataset is crucial, particularly when considering the further applications of IQA. As shown in Fig. 1(d), ATTIQA demonstrates outstanding generalization capabilities in cross-dataset validation experiments. Building on this strength, we propose a couple of applications of our IQA method. We show that our method can be used to evaluate the outputs of a generative model [33] and as a reward for reinforcement learning-based image enhancement [38].

## 2   Related Works

**Classical Image Quality Assessment.** Since image quality is highly regarded as essential in diverse vision applications, numerous image quality assessment studies have been explored. Traditional NR-IQA utilizes a feature-based machine learning approach to quantify image quality, leading to a primary focus on extracting meaningful features. Therefore, these works introduced hand-crafted feature based IQA, which is derived from natural scene statistics [8], spatial domain [25, 26] or frequency domain [35].

**Deep learning based IQA.** With the success of deep learning, various deep learning-based NR-IQA methods have been introduced. Early works tried to train convolutional neural networks by directly predicting mean opinion score (MOS) [15] or the distribution of MOS [41]. Some works attempted to incorporate meta-learning [53] or hypernetwork [40]. However, the limited size of IQA

datasets restricts these deep learning-based approaches, making it challenging to extract rich representations solely relying on the IQA dataset. To address this problem, recent IQA methods [41, 44, 50] commonly adopted ImageNet [34] classification-based backbone as their initial state, which already possesses rich representations. However, there is another problem that this representation is not fully suitable to IQA, as their pretraining task mainly focuses on semantic information but not image quality.

Beyond applying deep learning strategies to IQA, some approaches have focused on generating quality-aware representation using large datasets without the need for ground truth. Liu *et al.* [21] employed the Siamese network to rank images according to their quality, generating images of varying quality levels by applying different scales of distortion to a single image. Synthetic distortion-aware representation was introduced in [50], which attempts to classify the type or the amount of distortions applied to images. Recently, with the success of SSL, some works [23, 36, 52] suggested a contrastive learning framework refined for IQA. Unlike typical contrastive learning, they viewed patches from the same image, and each patch is differently distorted as different-quality samples. By treating these samples as negative pairs in the training process, they efficiently constructed a quality-aware representation space.

**Vision Language Model based IQA.** VLM [20,32] is a foundation model that learns correspondence between image and text to understand the relationships between visual contents and language. Specifically, CLIP [31] is trained with 400 million image-text pairs, and thus, it shows generalization capability for various computer vision tasks. Taking advantage of this ability, IQA methods that utilize CLIP have also been introduced. CLIP-IQA [43] directly assessed image quality by measuring the image's correspondence with quality-aware text prompts. They also suggested an enhanced version named CLIP-IQA+ that optimizes text prompts to adapt to the given target dataset. Despite this successful application, since the CLIP is trained with unrefined caption data focusing on semantic information, there is still room for improvement in refining CLIP's representation space towards an image quality-aware representation space. In response to this property, some works have tried to adapt the representation space of CLIP with additional datasets. Ke *et al.* [17] fine-tuned the CoCa [48] with the aesthetic captioning dataset to make aesthetic-aware VLM. By injecting aesthetic-related vision-language correspondence into VLM, they showed improved performance of their representation in quality-related downstream tasks. Zhang *et al.* [51] suggested a multi-task learning approach to adapt the vision backbone of VLM for a unified IQA dataset. They trained the vision backbone by optimizing cosine similarities of multi-modal embeddings with task-aware prompts.

As demonstrated by the above approaches, works refining the representation of VLM to aware image quality are currently underway. While these methodologies showed improvements in incorporating quality-aware information into VLM's representation space, the necessity of additional datasets to fine-tune remains a limitation. To mitigate this issue, our method does not *fine-tune* the
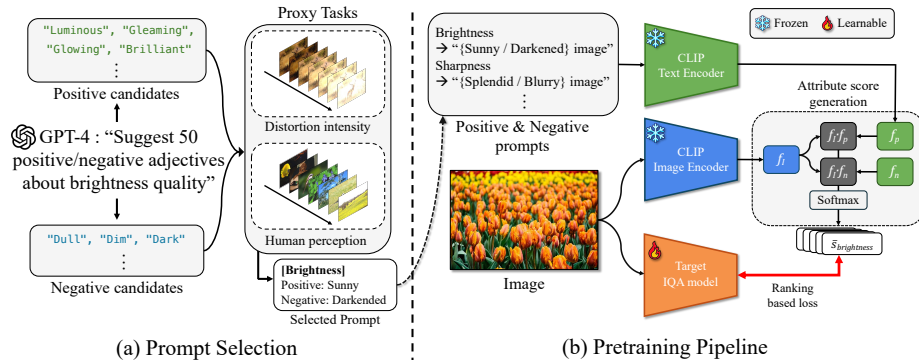
**Fig. 2:** (a) The overall process of our prompt selection strategy for each image attribute (e.g. brightness). Given attribute, we create prompt candidates using GPT-4 and then find the optimal prompt by utilizing proxy tasks related to the attribute. (b) ATTIQA's proposed pretraining pipeline. We generate attribute scores using CLIP with an antonym strategy and then train our target IQA model using ranking-based loss with generated scores.

CLIP model itself, but rather *extract* quality-related information from CLIP and use them to train our IQA model.

## 3   Methods

Fig. 2 illustrates our framework, which consists of two primary components: prompt selection and pretraining pipeline using pseudo-labels from CLIP [31]. During the prompt selection phase, we start by creating a list of candidate prompts using a large language model (LLM). We then identify the most suitable prompts for generating image attribute scores by evaluating the score generation ability of candidates by conducting proxy tasks. In the pretrain stage, we employ CLIP with chosen prompts to generate pseudo-labels for quantifying various image attribute scores on a large-scale dataset like ImageNet  [34]. Our IQA model is pretrained on this pseudo-labeled data and fine-tuned using a dedicated IQA dataset.

### 3.1   Image Attributes

Our method is designed based on CLIP guidance using specified prompts, enabling us to provide a new type of supervision beyond the Mean Opinion Score (MOS). Therefore, to reduce the ambiguity of the IQA task, which is typically represented solely by MOS, and to train our IQA model with more detailed criteria, we identify specific image attributes that affect the overall image quality as supervision. This approach aims to yield a more precise and well-defined representation of image quality by offering specific quality criteria beyond the generic and ambiguous mere notions of 'good' and 'bad'.

In line of IQA research that aims to incorporate quality relevant information beyond the MOS, five key attributes – *sharpness, contrast, brightness, color-fulness, and noisiness*– are widely employed and have proven beneficial for the IQA task [6,14,39]. Especially, SPAQ [6] substantiated through a user study that these five image attributes correlate well with perceived image quality. These observations indicate that these attributes are helpful factors in understanding the image quality. Therefore, to develop a detailed representation of IQA, we choose these five attributes as the target objective for our IQA model. Note that the following attribute score generation and prompt selection are conducted separately for each image attribute.

### 3.2    Attribute Score Generation

During attribute score generation, we generate five attribute scores for given images using CLIP's zero-shot inference. Given image $x$ and attribute-aware prompt $t$, CLIP encodes the image and prompt into shared multi-modal feature space. We then compute the relatedness score $s$ between $x$ and $t$ using cosine similarity as follows:

$$s(x,t) = \frac{E_I(x) \cdot E_T(t)^T}{||E_I(x)|| \cdot ||E_T(t)||},$$

(1)

where $E_I$ and $E_T$ represent the image encoder and text encoder of CLIP, respectively.

Our pseudo-label generation employs an antonym strategy [43], which computes scores by integrating scores of positive and negative prompts with the softmax function. For example, we can use a prompt pair {"*Dark image*", "*Bright image*"} as a negative-positive pair to calculate the brightness attribute score. Then, our attribute score is computed by the following equation:

$$\bar{s}_{attribute}(x) = \frac{e^{s(x,t_{pos})}}{e^{s(x,t_{pos})} + e^{s(x,t_{neg})}},$$

(2)

where $t_{pos}$ and $t_{neg}$ represents positive and negative prompt for the corresponding image attribute, respectively.

### 3.3    Prompt Selection

Previous work involving CLIP has shown that the choice of prompts is critical in determining performance. To address this, we introduce a prompt selection strategy aimed at identifying the most effective prompts for generating image attribute scores. Drawing inspiration from techniques used in the NLP field [7], we develop an efficient approach for identifying the optimal prompts.

As depicted in the left side of Figure 2(a), we begin by generating prompt candidates using Large Language Model (LLM), specifically GPT-4 [29]. To streamline the search process, we adopt a standard template for these prompts

in the format *"[adjective] image"*, focusing specifically on a variation of adjectives. Prompt candidates are constructed using GPT-4, with an ask query designed to elicit adjectives pertinent to specific image attributes. Since we use antonym pair for each attribute, we generate 50 positive and 50 negative adjectives, resulting in 2500 positive-negative prompt candidates for each attribute.

Subsequently, we identify the most suitable prompt pair from 2500 candidates by assessing their capability to generate accurate attribute scores. To find the optimal prompt pair, we present two proxy tasks. The optimal prompt is determined as the prompt that produces an attribute score that best aligns with the goals of both tasks. To measure an image attribute appropriately, the proxy tasks are de-

**Table 1:** Results of the prompt selection. These prompts are chosen by our prompt selection strategy.

| Attribute | Positive prompt | Negative prompt |
|---|---|---|
| Sharpness | *"Splendid image"* | *"Blurry image"* |
| Contrast | *"Distinct image"* | *"Vague image"* |
| Brightness | *"Sunny image"* | *"Darkened image"* |
| Colorfulness | *"Vibrant image"* | *"Colorless image"* |
| Noisiness | *"Peerless image"* | *"Grainy image"* |

signed under two hypotheses: 1) For a fixed image, when a distortion corresponding to the image attribute is applied to it, the attribute score predicted by the prompt should increase or decrease accordingly. 2) For different images, the attribute scores generated by the prompt pairs should match well with the degree of human perception of each attribute.

In the first proxy task, we apply varying levels of distortion to a fixed image and identify the optimal pair of prompts that yield an attribute score aligning most accurately with the applied level of distortion. For the second task, we employ the SPAQ dataset, which comprises diverse scenes and offers human-annotated attribute scores for the same five attributes we adopt. We aim to identify the prompt pair that generates the attribute scores whose order closely aligns with the order of the provided scores in the dataset. We calculate the sum of SROCC scores for both tasks and select the highest performing prompt pair, and the result is shown in Table 1. These selected pairs are then utilized to generate each attribute score in the following pretraining pipeline.

### 3.4   Pretraining Pipeline

After the prompt selection, we train the target IQA model to construct attribute-aware space with ranking-based loss using a pseudo-label derived from CLIP, as illustrated in Fig. 2(b).

Our method seeks to create five unique representation spaces for each specific image attribute, so we adapt the architecture of the target IQA model to incorporate these characteristics. Accordingly, our IQA model comprises a shared backbone and five attribute heads for each image attribute. Each attribute head consists of two-layer multi-layer perceptrons (MLPs) that aim to output an attribute score. Then, our training objective for the pretraining pipeline can be formulated by minimizing the discrepancy between five attribute score predictions from the IQA model and the corresponding image attribute scores generated from CLIP.

Our pretraining pipeline can be directly implemented using regression-based loss such as MSE or L1 loss. However, there are some problems with directly using the attribute score with regression-based loss. Assessing image attributes in a zero-shot manner is an inherently challenging task, and the representation space of CLIP is not specifically tailored for this task. Due to these problems, our attribute score can be regarded as a noisy label. Therefore, regression-based loss may hinder the pretraining pipeline and result in diminished performance. Therefore, we use a relative ranking-based loss instead of a numerical norm-based loss to minimize the dependence on CLIP's numerical results, which are subject to uncertainty. To implement this loss in our framework, we utilize margin-ranking loss that optimizes the relative ranking of the two samples given. We first define indicator function $F$, which specifies the superiority of image attribute based on its score $\bar{s}$ for given images:

$$F_a(x_1, x_2) = \begin{cases} 0, & \bar{s}_a(x_1) > \bar{s}_a(x_2) \\ 1, & \bar{s}_a(x_1) \le \bar{s}_a(x_2), \end{cases} \tag{3}$$

where $a$ denotes an element of image attributes set $A = \{sharpness,\ contrast,\ brightness,\ colorfulness,\ noisiness\}$, and $x_1$ and $x_2$ denote the sample images.

Then, we train our target IQA model based on margin-ranking loss with the indicator function $F_a$. We compute our loss by summation of margin-ranking loss independently for each attribute $a$ using attribute score prediction from respective attribute head $E_a$:

$$L = \sum_a \max(0, m - (E_a(x_1) - E_a(x_2))) \cdot F_a(x_1, x_2), \tag{4}$$

where $m$ denotes the margin hyper-parameter.

### 3.5 Fine-tuning

To predict MOS with our IQA model, we have to fine-tune it on the target IQA dataset. However, the architecture of our IQA model during the pretraining stage is not designed to output a single score but instead predicts five attribute scores. Therefore, we adapt the architecture of our model to fit the IQA task to generate a single MOS as output. At a fine-tuning stage, we extract features from each attribute head, excluding the final layer, and these features are then concatenated and fed into regression MLPs that predict the MOS.

## 4  Experiments

### 4.1  Datasets

We conduct experiments with ATTIQA on the 4 "in-the-wild" NR-IQA Datasets, CLIVE [9], KonIQ-10k [13], SPAQ [6], and FLIVE [46] and 1 image aesthetic dataset, AVA [27].

CLIVE includes 1,162 authentically distorted images collected from real-world mobile camera devices. KonIQ-10k [13] consists of 10,073 images sourced

from dataset YFCC100m [42]. This dataset focuses on the intrinsic quality of the image by cropping around meaningful areas and comprises a uniform sampling process that considers image tags, resolution, and quality-related statistics. SPAQ [6] comprises 11,125 images captured with multiple cameras. They conducted a user study focused on five pre-defined image attributes, providing numerical scores for each attribute with the MOS of given images. FLIVE [46] consists of 40,000 images and 120,000 patches cropped from each image. These images are sourced from AVA [27], VOC [5], EMOTIC [18], and CERTH Blur [24]. Note that our experiment only utilizes 40,000 original images, not using patches. AVA [27] is an image aesthetic dataset comprising 250,000 images. While this dataset aims to measure image aesthetics, we utilize the AVA [27] for our experiments since its user study setting is the same as "IQA in the wild", and there is an intrinsically ambiguous boundary between quality and aesthetics.

For FLIVE [46] and AVA [27], we follow the official dataset split. For the rest, we randomly partition the dataset and allocate 80% to the train set and 20% to the test set. To compensate for the bias arising from random splits, following previous works, we conduct the same experiment ten times with different random splits and report the median value as a result.

### 4.2  Experimental Setup

**Implementation detail.** For a fair comparison with baselines, we utilize ResNet-50 [11] as our backbone, widely used in pretraining-based IQA. The attribute head is a two-layer MLPs whose hidden channel is 512. For CLIP, we adopt the ViT-B/16 model [4]. For the margin parameter $m$ at the loss function, We experimentally set the value of $m$ as 0.1.

**Pretraining.** Since our method does not require any ground truth at this stage, we use ImageNet dataset [34] widely used for pretext tasks. We resized the image's shorter edge to 256 and randomly cropped the image at a resolution of $224 \times 224$. In this process, we use AdamW optimizer [22] with weight decay 0.01 and set batch size 256. We train our network for 100 epochs and set an initial learning rate of 1e-4, which decays 0.1 scale at 60 and 80 epochs.

**Fine-tuning.** At the fine-tuning stage, we followed the setting from [52]. We resized the image's shorter edge to 340 and randomly cropped the image at a resolution of $320 \times 320$. We use AdamW optimizer [22] with weight decay 0.01 and set batch size 64. We fine-tune our network 100 epochs and set initial learning rates 1e-4, 5e-5, and 1e-5 on CLIVE, KonIQ, and the rest, respectively, with the cosine annealing decay strategy.

**Evaluation.** At the evaluation stage, we take five crops at a resolution of $320 \times 320$ from each corner and center as test samples, and the average of the results is used for the predicted MOS. For performance evaluation, we calculated Pearson's Linear Correlation Coefficient (PLCC) and Spearman's Rank-Order Correlation Coefficient (SROCC), which are widely adopted evaluation metrics in IQA research.

**Table 2:** fine-tuning performance comparison of ATTIQA and existing NR-IQA methods for 4 IQA "in-the-wild" dataset and 1 IAA dataset. "†" denotes that this measurement is achieved from [52]. "-" denotes that measurement is not possible due to the absence of an official code and result. Other measurements are based on the official reports or reproduced by the official code. We highlight the best performance in bold and underline the second-best performance for each dataset.

| Methods | CLIVE | | KonIQ | | SPAQ | | FLIVE | | AVA | |
|---|---|---|---|---|---|---|---|---|---|---|
| | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC |
| DBCNN† [50] | 0.844 | 0.862 | 0.878 | 0.887 | 0.906 | 0.907 | 0.542 | 0.626 | 0.554 | 0.583 |
| HyperIQA† [40] | 0.855 | 0.871 | 0.908 | 0.921 | 0.916 | 0.919 | 0.535 | 0.623 | 0.668 | 0.668 |
| CONTRIQUE [23] | 0.824 | 0.848 | 0.900 | 0.915 | 0.910 | 0.915 | 0.598 | 0.674 | 0.674 | 0.678 |
| MUSIQ [16] | - | - | 0.916 | 0.928 | 0.917 | 0.921 | **0.646** | 0.739 | 0.726 | 0.738 |
| TReS [10] | 0.846 | 0.877 | 0.915 | 0.928 | 0.915 | 0.919 | 0.554 | 0.625 | 0.658 | 0.663 |
| Re-IQA [36] | 0.823 | 0.865 | 0.924 | 0.935 | 0.915 | 0.919 | 0.574 | 0.674 | 0.714 | 0.716 |
| QPT [52] | <u>0.895</u> | 0.914 | 0.927 | 0.941 | <u>0.925</u> | <u>0.928</u> | 0.610 | 0.677 | - | - |
| Imagenet Pretrained [34] | 0.816 | 0.861 | 0.906 | 0.925 | 0.911 | 0.915 | 0.553 | 0.666 | 0.696 | 0.698 |
| CLIP [31] | 0.847 | 0.881 | 0.918 | 0.932 | 0.918 | 0.922 | 0.563 | 0.628 | 0.746 | 0.745 |
| CLIP-IQA+ [23] | 0.805 | 0.832 | 0.895 | 0.909 | 0.864 | 0.866 | 0.575 | 0.593 | 0.692 | 0.732 |
| LIQE [51] | 0.865 | 0.866 | 0.898 | 0.913 | - | - | - | - | - | - |
| ATTIQA (Distortion intensity) | 0.891 | 0.910 | <u>0.929</u> | <u>0.943</u> | 0.922 | 0.926 | 0.625 | 0.729 | 0.754 | 0.750 |
| ATTIQA (Human perception) | 0.890 | <u>0.915</u> | **0.942** | **0.952** | <u>0.925</u> | **0.930** | <u>0.635</u> | <u>0.740</u> | <u>0.756</u> | <u>0.759</u> |
| ATTIQA (Joint strategy) | **0.898** | **0.916** | **0.942** | **0.952** | **0.926** | **0.930** | 0.632 | **0.742** | **0.761** | **0.761** |

### 4.3 Main Result

In this section, we report the quantitative performance of ATTIQA and compare it with existing NR-IQA models. Utilizing our CLIP-guided attribute-aware pretrained model, we conduct fine-tuning on five IQA datasets to predict the MOS. As shown in the Table 2, ATTIQA shows notable performance improvements on four IQA "in the wild" datasets and one image aesthetic dataset AVA. Our method demonstrates state-of-the-art performance in most evaluation settings, with the second-best SROCC performance on the FLIVE dataset. It is important to note that MUSIQ is an exceptional work that utilized the complete FLIVE dataset comprising patches and full images, unlike the other methods that do not use patch data. Notably, ATTIQA shows a significant performance gap on the KonIQ-10k and AVA datasets.

In the last three rows of Table 2, we report the performance of three different versions of ATTIQA. We carry out experiments with various types of prompts, including cases where we apply the two proxy tasks described in Sec 3.3—Distortion Intensity and Human Perception—separately, as well as a scenario where we combine both tasks in our prompt selection strategy (Joint Strategy). We observe that prompts based on *human perception* work effectively, and the *joint strategy* that involves both proxy tasks shows the most superior performance. It indicates that a prompt selection strategy that considers using both low-level information *distortion* and high-level *human perception* enhances the robustness of our model across various datasets.

### 4.4 Generalization Capability

**Cross-dataset Validation.** To verify ATTIQA's generalization ability, we conduct experiments about cross-dataset validation, where the training and testing

**Table 3:** Cross dataset validation performance comparison of ATTIQA and existing NR-IQA methods.

| Train DB | CLIVE | | | KonIQ | | | SPAQ | | | FLIVE | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Test DB | KonIQ | SPAQ | FLIVE | CLIVE | SPAQ | FLIVE | CLIVE | KonIQ | FLIVE | CLIVE | KonIQ | SPAQ |
| CONTRIQUE | 0.676 | 0.842 | 0.346 | 0.731 | 0.789 | 0.410 | 0.549 | 0.646 | 0.338 | 0.706 | 0.709 | 0.734 |
| Re-IQA | 0.769 | 0.852 | 0.424 | 0.791 | 0.862 | 0.461 | 0.732 | 0.707 | 0.497 | 0.720 | 0.676 | 0.793 |
| CLIP | 0.725 | 0.850 | 0.405 | 0.799 | 0.837 | 0.507 | 0.773 | <u>0.753</u> | 0.496 | 0.727 | 0.717 | 0.834 |
| CLIP-IQA+ | 0.697 | 0.836 | 0.437 | 0.803 | 0.832 | 0.516 | 0.784 | 0.722 | 0.470 | 0.620 | 0.631 | 0.661 |
| LIQE | <u>0.819</u> | <u>0.877</u> | <u>0.497</u> | <u>0.824</u> | <u>0.868</u> | **0.551** | - | - | - | - | - | - |
| ATTIQA | **0.829** | **0.887** | **0.511** | **0.856** | **0.879** | <u>0.540</u> | **0.824** | **0.819** | **0.548** | **0.756** | **0.762** | **0.867** |

datasets are distinct. This experiment allows us to evaluate the IQA model's ability to learn generalizable features by training it on the specified dataset and testing it on the unseen dataset. To consider various scenarios, we conduct extensive experiments across four datasets: CLIVE, KonIQ, SPAQ, and FLIVE. Every experimental setup is the same as the main experiment, and due to the various ranges of the MOS for each dataset, we use only SROCC as an evaluation criterion. As shown in Table 3, ATTIQA exhibits superior generalization capability compared to baselines, achieving the best performance in most scenarios. Interestingly, LIQE achieves comparable results to ATTIQA, demonstrating that strategies adapting CLIP possess strong generalization capabilities. However, we highlight that ATTIQA outperforms LIQE in most scenarios without requiring additional annotations and while employing a more lightweight model architecture. This difference exhibits our method's enhanced suitability for real-world scenarios, indicating improved generalizability and scalability.

For real-world applications, a model's generalization ability is far more critical than its performance on specific benchmark datasets. Our method's superior generalization capability ensures robust baseline performance on unseen images, highlighting its practicability when considering the purpose of the IQA tasks. Building on this strength, we will demonstrate the application of ATTIQA in real-world scenarios in Sec 5.

### 4.5   Ablation Studies

**Linear probing.** We conduct linear probing experiments to demonstrate the robustness of our attribute-aware pretrained feature space, training only a single regression MLPs on the target dataset with a frozen ATTIQA backbone. In this experiment, we compared ATTIQA to models that mainly focus on learning the representation space for IQA. For Re-IQA, we report the three types of results based on the encoder configuration: using only the quality encoder, the content encoder, and both encoders. As shown in Table 4, ATTIQA shows a significant performance gap compared to other methods in CLIVE. In other datasets, ATTIQA also demonstrates comparable results to Re-IQA, which uses both features of a separate quality and content encoder, while our method only utilizes a single shared encoder for five attributes.

**Table 4:** Linear probing performance comparison of ATTIQA and NR-IQA methods which focuses on representation learning.

| Methods | CLIVE | | KonIQ | | SPAQ | |
|---|---|---|---|---|---|---|
| | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC |
| CONTRIQUE | 0.845 | 0.857 | 0.894 | 0.906 | 0.916 | 0.919 |
| Re-IQA (quality) | 0.806 | 0.824 | 0.861 | 0.885 | 0.900 | 0.910 |
| Re-IQA (content) | 0.808 | 0.844 | 0.896 | 0.912 | 0.902 | 0.908 |
| Re-IQA (both) | 0.840 | 0.854 | **0.914** | **0.923** | **0.918** | **0.925** |
| CLIP | 0.803 | 0.829 | 0.883 | 0.895 | 0.895 | 0.896 |
| ATTIQA | **0.870** | **0.891** | 0.903 | 0.918 | 0.918 | 0.922 |

**Table 5:** Ablation study results about our prompt based strategy and loss function.

| Prompt type | CLIVE | | KonIQ | | SPAQ | |
|---|---|---|---|---|---|---|
| | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC |
| ATTIQA | **0.898** | **0.916** | **0.942** | **0.952** | **0.926** | **0.930** |
| single-prompt | 0.880 | 0.909 | 0.928 | 0.939 | 0.916 | 0.920 |
| Worst prompt | 0.869 | 0.889 | 0.930 | 0.943 | 0.920 | 0.925 |
| Median prompt | 0.872 | 0.893 | 0.931 | 0.944 | 0.921 | 0.925 |
| with $L_2$ | 0.875 | 0.904 | 0.933 | 0.945 | 0.923 | 0.928 |

**Attribute based Approach.** To demonstrate the effectiveness of our image attribute based approach, we carry out an experiment by replacing the target objective from five image attributes with a single overall image quality. In this experiment, we train the IQA model with a single pseudo-label using a prompt pair that describes image quality: *{"Good image", "Bad image"}*. Comparing the first and the second tab of Table 5, we can observe that our representation space decomposing image quality into five attributes outperforms the single-prompt based representation space, justifying our approach for model design.

Moreover, we conduct an additional experiment to verify that each attribute head accurately captures its corresponding attributes. The experiment uses Grad-CAM [37] to explore how the importance of each attribute head used for MOS prediction varies when only one attribute is modified for a given image. The experiment shows that if a particular attribute varies excessively to the extent that compromises the quality of the image, the head of that attribute has a more significant impact on MOS estimation. This analysis indicates that each attribute head precisely captures individual attributes and the response to each attribute variation is effectively reflected in the MOS prediction. For detailed experimental setup and visualized results, please refer to the supplementary materials.

**Prompt Selection Strategy.** To justify our prompt selection strategy, we conduct experiments with other prompts selected by different strategies. For comparison, we adopt prompts that achieve the median and lowest scores in the proxy task. As shown in the third tab of Table 5, the results of our strategy align with the performance at the evaluation. This correlation validates the efficacy of our prompt selection strategy.

**Ranking-based loss.** To verify the efficacy of our relative ranking-based loss approach, we conduct an additional ablation study by replacing the margin-ranking loss with L2 loss at the pretraining stage. As depicted in the last row of Table 5, the use of L2 loss exhibits a performance degradation compared to adopting margin-ranking loss. Interestingly, we can observe a notable performance decline in the CLIVE dataset, which has the smallest dataset size within this ablation study. This result supports the use of relative loss instead of numerical loss, enhancing our pretraining pipeline's robustness.

**Table 6:** Comparisons of accuracy between human preferences and IQA model's result.

| Method | CONTRIQUE | Re-IQA | CLIP-IQA+ | ATTIQA |
|---|---|---|---|---|
| Accuracy (%) | 61.5 | 55.0 | 57.5 | **71.0** |

**Table 7:** Performance comparisons among IQA model in an AI-Generated Contents Dataset(AGIQA-3k)

| Method | CONTRIQUE | Re-IQA | CLIP-IQA+ | ATTIQA |
|---|---|---|---|---|
| SROCC | 0.643 | 0.807 | 0.835 | **0.854** |
| PLCC | 0.795 | 0.876 | 0.885 | **0.911** |



**Fig. 3:** Example of generated images. The images are generated by the same prompt. ATTIQA hits human preference while others do not.

## 5    Applications

To better demonstrate our ATTIQA's generalization capability, we introduce two types of applications in this section: (1) metrics for the generative model and (2) image enhancement guided by our IQA score. For each application, we employ models trained on the KonIQ dataset, which shows the best generalization capability at Sec 4.4.

### 5.1    Metrics for Generative Model

Recently, as the diffusion models [12,33] have shown success in the image generation task, there is a growing body of research in text-to-image generation [28,32]. One of their primary focus is generating high quality images from a given identical text prompt. In this regard, we attempt to employ ATTIQA as a metric for generative models.

To assess ATTIQA's effectiveness as a metric, we create a benchmark dataset that involves the pairwise comparison of two images generated from the same text prompt. Here, we generate 200 pairs of images using the Stable Diffusion [33], widely used in text-to-image synthesis, and collect human preference by conducting a user study. When collecting the user preferences, we only present the generated images without the prompt to make participants focus on visual quality. The user study was carried out with 60 participants through Amazon Mechanical Turk (AMT). We then investigate the correlation between IQA models and the human participants.

As shown in Table 6, our method mostly aligns with human preference compared to other IQA methods. Images A and B are synthesized using the same prompt, but their visual quality varies from image to image. Our ATTIQA can capture this detailed visual quality difference while others do not. Please refer to the supplementary for the user study details and more visual results. We will make the benchmark used in this application publicly available for further IQA research.

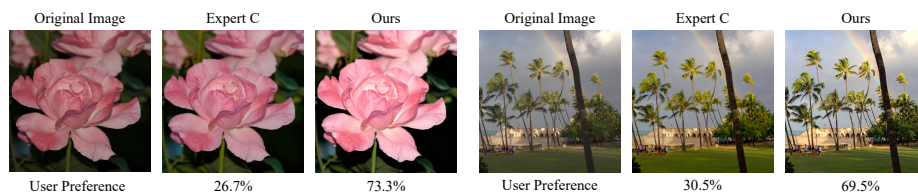| Original Image | Expert C | Ours | Original Image | Expert C | Ours |
| --- | --- | --- | --- | --- | --- |
| User Preference | 26.7% | 73.3% | User Preference | 30.5% | 69.5% |

**Fig. 4:** Qualitative comparisons between our image enhancement model and retouching of Expert C. Our result gives more liveliness and vibrancy, aligned more closely with human preference.

Moreover, we carry out an additional experiment using an AGIQA-3k dataset [19], which consists of images generated by various generative models. As shown in Table 7, ATTIQA outperforms other methods, exhibiting a significant performance gap. These results highlight the improved generalization capability of our method when extended to AI-generated content. They indicate the potential for expanding the use of ATTIQA as a metric to evaluate generative models.

### 5.2   Image Enhancement

The image signal processing pipeline (ISP) converts an input raw image into a color image. It is essential to carefully tune the parameters of the ISP to obtain visually pleasing images. In this section, we apply ATTIQA's MOS prediction as a reward for reinforcement learning to find optimal parameters for the ISP [38]. After the training, we convert raw images into color images in the MIT-Adobe-5k dataset [3], which consists of 5,000 raw images and color images retouched by five experts (A/B/C/D/E). Then, we conduct a user study comparing our result against the retouched one by expert C, which is typically used as the ground truth in most previous image enhancement research. The study was executed with 60 participants through AMT, involving a comparison of 200 image pairs. For details on the implementation, please refer to the supplementary materials.

As shown in Fig. 4, our pipeline retouches images to make them more colorful and vivid compared to both retouching by expert C and the default settings. Furthermore, according to our user study, ATTIQA receives higher preferences from subjects, demonstrating a 58% win rate compared to Expert C. We also report additional qualitative comparisons to supplementary material.

## 6   Discussion and Conclusion

We propose ATTIQA, a pretraining framework for IQA that develops an attribute-aware representation space with CLIP guidance. Since our IQA model effectively incorporates CLIP's vast knowledge and scalability of large datasets, it shows state-of-the-art performance on IQA datasets and superior generalization capability on cross-dataset validation. Leveraging these advantages, we successfully demonstrate a couple of real-world applications where IQA can be utilized.

**Limitation and Future Work.** While our approach focuses on five attributes commonly employed in the IQA domain, we may expect that other properties relevant to image quality exist (e.g., Composition and Focus). Consequently, future work will involve exploring extended representation spaces for IQA. Given that our method proposes a pretraining framework not limited to the specified attributes, our work holds the potential for expansion to encompass additional properties.

# References

1. Bianco, S., Celona, L., Napoletano, P., Schettini, R.: On the use of deep learning for blind image quality assessment. Signal, Image and Video Processing **12**, 355–362 (2018) 2
2. Bosse, S., Maniry, D., Müller, K.R., Wiegand, T., Samek, W.: Deep neural networks for no-reference and full-reference image quality assessment. IEEE Transactions on Image Processing (TIP) **27**(1), 206–219 (2017) 1
3. Bychkovsky, V., Paris, S., Chan, E., Durand, F.: Learning photographic global tonal adjustment with a database of input / output image pairs. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2011) 14
4. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. In: International Conference on Learning Representations (ICLR) (2021) 9
5. Everingham, M., Eslami, S.M.A., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A.: The pascal visual object classes challenge: A retrospective. International Journal of Computer Vision (IJCV) **111**, 98–136 (2015) 9
6. 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 (CVPR). pp. 3677–3686 (2020) 1, 6, 8, 9
7. Gao, T., Fisch, A., Chen, D.: Making pre-trained language models better few-shot learners. In: Zong, C., Xia, F., Li, W., Navigli, R. (eds.) Proceedings of the Association for Computational Linguistics (ACL). pp. 3816–3830 (2021) 6
8. Gao, X., Gao, F., Tao, D., Li, X.: Universal blind image quality assessment metrics via natural scene statistics and multiple kernel learning. IEEE Transactions on neural networks and learning systems (2013) 3
9. Ghadiyaram, D., Bovik, A.C.: Massive online crowdsourced study of subjective and objective picture quality. IEEE Transactions on Image Processing (TIP) **25**(1), 372–387 (2015) 1, 8
10. 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 (WACV). pp. 1220–1230 (2022) 1, 10
11. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 770–778 (2016) 9
12. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems (NeurIPS) **33**, 6840–6851 (2020) 13

13. 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 (TIP) **29**, 4041–4056 (2020) 1, 8

14. Huang, Y., Li, L., Yang, Y., Li, Y., Guo, Y.: Explainable and generalizable blind image quality assessment via semantic attribute reasoning. IEEE Transactions on Multimedia (TMM) (2022) 6

15. Kang, L., Ye, P., Li, Y., Doermann, D.: Convolutional neural networks for no-reference image quality assessment. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1733–1740 (2014) 1, 3

16. Ke, J., Wang, Q., Wang, Y., Milanfar, P., Yang, F.: Musiq: Multi-scale image quality transformer. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 5148–5157 (2021) 10

17. Ke, J., Ye, K., Yu, J., Wu, Y., Milanfar, P., Yang, F.: Vila: Learning image aesthetics from user comments with vision-language pretraining. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 10041–10051 (2023) 4

18. Kosti, R., Alvarez, J., Recasens, A., Lapedriza, A.: Context based emotion recognition using emotic dataset. IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI) (2019) 9

19. Li, C., Zhang, Z., Wu, H., Sun, W., Min, X., Liu, X., Zhai, G., Lin, W.: Agiqa-3k: An open database for ai-generated image quality assessment (2023) 14

20. Li, J., Li, D., Xiong, C., Hoi, S.: Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In: International conference on machine learning (ICML). pp. 12888–12900 (2022) 4

21. Liu, X., Van De Weijer, J., Bagdanov, A.D.: Rankiqa: Learning from rankings for no-reference image quality assessment. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 1040–1049 (2017) 1, 2, 4

22. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: International Conference on Learning Representations (ICLR) (2018) 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 (TIP) **31**, 4149–4161 (2022) 2, 4, 10

24. Mavridaki, E., Mezaris, V.: No-reference blur assessment in natural images using fourier transform and spatial pyramids. In: IEEE International Conference on Image Processing (ICIP). pp. 566–570 (2014) 9

25. Mittal, A., Moorthy, A.K., Bovik, A.C.: No-reference image quality assessment in the spatial domain. IEEE Transactions on Image Processing (TIP) **21**(12), 4695–4708 (2012) 1, 3

26. Mittal, A., Soundararajan, R., Bovik, A.C.: Making a "completely blind" image quality analyzer. IEEE Signal Processing Letters **20**(3), 209–212 (2012) 3

27. Murray, N., Marchesotti, L., Perronnin, F.: Ava: A large-scale database for aesthetic visual analysis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2408–2415 (2012) 8, 9

28. Nichol, A., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., McGrew, B., Sutskever, I., Chen, M.: Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In: International conference on machine learning (ICML) (2021) 13

29. OpenAI: Gpt-4 technical report (2023) 6

30. Ponomarenko, N.N., Jin, L., Ieremeiev, O., Lukin, V.V., Egiazarian, K.O., Astola, J., Vozel, B., Chehdi, K., Carli, M., Battisti, F., Kuo, C.C.J.: Image database tid2013: Peculiarities, results and perspectives. Signal Processing, Image Communication **30**, 57–77 (2015) 1

31. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. In: International conference on machine learning (ICML) (2021) 2, 4, 5, 10

32. Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., Chen, M., Sutskever, I.: Zero-shot text-to-image generation. In: International conference on machine learning (ICML). pp. 8821–8831 (2021) 4, 13

33. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 10684–10695 (2022) 3, 13

34. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M.S., Berg, A.C., Fei-Fei, L.: Imagenet large scale visual recognition challenge. International Journal of Computer Vision (IJCV) pp. 211 – 252 (2014) 2, 4, 5, 9, 10

35. 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 (TIP) **21**(8), 3339–3352 (2012) 3

36. 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 (CVPR). pp. 5846–5855 (2023) 2, 4, 10

37. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Gradcam: Visual explanations from deep networks via gradient-based localization. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 618–626 (2017) 12

38. Shin, U., Lee, K., Kweon, I.S.: Drl-isp: Multi-objective camera isp with deep reinforcement learning. In: IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). pp. 7044–7051 (2022) 3, 14

39. Su, S., Hosu, V., Lin, H., Zhang, Y., Saupe, D.: Koniq++: Boosting no-reference image quality assessment in the wild by jointly predicting image quality and defects. In: British Machine Vision Conference (BMVC) (2021) 6

40. 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) 1, 3, 10

41. Talebi, H., Milanfar, P.: Nima: Neural image assessment. IEEE Transactions on Image Processing (TIP) **27**(8), 3998–4011 (2018) 1, 2, 3, 4

42. Thomee, B., Shamma, D.A., Friedland, G., Elizalde, B., Ni, K., Poland, D., Borth, D., Li, L.J.: Yfcc100m: The new data in multimedia research. Communications of the ACM **59**(2), 64–73 (2016) 9

43. 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 (AAAI). vol. 37, pp. 2555–2563 (2023) 2, 3, 4, 6

44. Yang, X., Li, F., Liu, H.: Ttl-iqa: Transitive transfer learning based no-reference image quality assessment. IEEE Transactions on Multimedia (TMM) **23**, 4326–4340 (2021). https://doi.org/10.1109/TMM.2020.3040529 2, 4

45. Ye, P., Kumar, J., Kang, L., Doermann, D.: Unsupervised feature learning framework for no-reference image quality assessment. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1098–1105 (2012) 1

46. Ying, Z., Niu, H., Gupta, P., Mahajan, D., Ghadiyaram, D., Bovik, A.: From patches to pictures (paq-2-piq): Mapping the perceptual space of picture quality. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3575–3585 (2020) 1, 8, 9

47. You, J., Korhonen, J.: Transformer for image quality assessment. In: IEEE International Conference on Image Processing (ICIP). pp. 1389–1393 (2021) 1

48. Yu, J., Wang, Z., Vasudevan, V., Yeung, L., Seyedhosseini, M., Wu, Y.: Coca: Contrastive captioners are image-text foundation models (2022) 4

49. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 586–595 (2018) 1

50. Zhang, W., Ma, K., Yan, J., Deng, D., Wang, Z.: Blind image quality assessment using a deep bilinear convolutional neural network. IEEE Transactions on Circuits and Systems for Video Technology (TCVST) **30**(1), 36–47 (2018) 1, 2, 4, 10

51. 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 (CVPR). pp. 14071–14081 (2023) 2, 4, 10

52. Zhao, K., Yuan, K., Sun, M., Li, M., Wen, X.: Quality-aware pre-trained models for blind image quality assessment. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 22302–22313 (2023) 2, 4, 9, 10

53. Zhu, H., Li, L., Wu, J., Dong, W., Shi, G.: Metaiqa: Deep meta-learning for no-reference image quality assessment. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 14143–14152 (2020) 1, 3

# (Supplementary material)
# CLIP-Guided Attribute Aware Pretraining for Generalizable Image Quality Assessment

Daekyu Kwon[1], Dongyoung Kim[1], Sehwan Ki[2], Younghyun Jo[2],
Hyong-Euk Lee[2], and Seon Joo Kim[1]

[1] Yonsei University
[2] Samsung Advanced Institue of Technology

For supplementary materials, we provide societal impacts, details on our implementation and additional experiment results.

## A    Societal Impacts

Our work contributes to enhancing the quality of digital content by introducing a powerful tool for evaluation across diverse platforms, such as social media, online publishing, and digital marketing. Furthermore, we propose a robust assessment tool that can employed across different image domains, enabling application in fields like medical imaging, remote sensing, and surveillance. By elevating the quality and reliability of images for these fields, they can deliver more precise analyses and decisions, potentially saving lives and conserving resources.

While enhanced IQA has many positive applications, it could also, conversely, help generate more offensive content, such as deepfakes. This work might be misused for creating false narratives or misinformation, impacting public opinion and individual reputations. Moreover, developing advanced image quality assessment models may require significant computational resources. This can lead to the digital divide, where only well-resourced organizations can afford to exploit this technology and raise concerns about the environmental impact of increased energy consumption.

## B    Details on Prompt Selection

This section explains the prompts and techniques used in the prompt selection strategy. Our prompt selection strategy generates prompt candidates from GPT-4 [?] through the following query that receives a set of adjectives: *"Suggest 50 positive/negative adjectives about {attribute} related to image quality"*. We create prompt candidates by changing *{attribute}* in the query to aligned attribute. The generated prompt candidates are reported in Table 13.

We also provide details on the proxy tasks stage. For the distortion intensity based proxy task, which measures attribute score by predicting the intensity of the corresponding distortion, we pair image attributes and corresponding distortions as follows:

**Table 8:** Results of the prompt selection with single proxy task. We only denote adjective for this table.

| Proxy task | Distortion intensity | | Human perception | |
|---|---|---|---|---|
| Attribute | Positive | Negative | Positive | Negative |
| Sharpness | *"Unambiguous"* | *"Vague"* | *"High-definition"* | *"Out-of-focus"* |
| Contrast | *"Enhanced"* | *"Bleak"* | *"Splendid"* | *"Blurred"* |
| Brightness | *"Clear"* | *"Starless"* | *"Dark"* | *"High-key"* |
| Colorfulness | *"Multicolored"* | *"Grayish"* | *"Lively"* | *"Blurred"* |
| Noisiness | *"Clutter-free"* | *"Spattered"* | *"Peerless"* | *"Blurry"* |

- Sharpness : Gaussian Blur, ZoomBlur, LensBlur
- Contrast : Contrast adjustment multiplying factors to RGB value
- Brightness : V adjustment in HSV space
- Colorfulness : Saturation adjustment in HSV space
- Noisiness : Gaussian Noise, ISO Noise

In attributes where multiple distortions are described, we execute proxy tasks for each distortion and compute the final result by averaging each output. In Table 8, we note selected prompts that only utilize a single proxy task whose results are reported in Table 2.

## C   Additional Experiments

### C.1   Impact of dataset scale

To verify the scalability of our pretraining scheme, we conduct experiments changing the ratio of the used dataset. In this experiment, we pretrain our model only utilizing 20% and 50% of the ImageNet dataset. As shown in Table 9, the performance of ATTIQA increases with the amount of available datasets. This result indicates that the performance of ATTIQA can be improved when we can train it on larger datasets (e.g., ALIGN, JFT-300M).

**Table 9:** Performance comparison of ATTIQA using various ratios of datasets.

| Methods | CLIVE [?] | | KonIQ [?] | | SPAQ [?] | |
|---|---|---|---|---|---|---|
| | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC |
| 20% | 0.887 | 0.904 | 0.935 | 0.946 | 0.923 | 0.928 |
| 50% | 0.889 | 0.913 | 0.938 | 0.949 | 0.924 | 0.928 |
| 100% | **0.898** | **0.916** | **0.942** | **0.952** | **0.926** | **0.930** |

**Table 10:** performance comparison of fine-tuning which only utilizes ATTIQA's each attribute head.

| Methods | CLIVE | | KonIQ | | SPAQ | |
|---|---|---|---|---|---|---|
| | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC |
| Sharpness | 0.890 | 0.910 | 0.939 | 0.950 | 0.925 | 0.929 |
| Contrast | 0.870 | 0.905 | 0.937 | 0.947 | 0.923 | 0.928 |
| Brightness | 0.880 | 0.905 | 0.936 | 0.948 | 0.923 | 0.926 |
| Colorfulness | 0.881 | 0.906 | 0.934 | 0.945 | 0.924 | 0.927 |
| Noisiness | 0.878 | 0.901 | 0.939 | 0.950 | 0.924 | 0.928 |
| ATTIQA | **0.898** | **0.916** | **0.942** | **0.952** | **0.926** | **0.930** |

### C.2    Impact of attribute head

To evaluate the impact of each attribute head, we conduct experiments that only utilize the head's feature space. To implement this experiment, we only adopt a single attribute head for fine-tuning instead of concatenating the five attributes feature map. As shown in Table 10, although each attribute head shows similar performance since they share a backbone, we observe a slightly better performance in the sharpness attribute head. In contrast, the performance of other attribute heads varies depending on the dataset. It can be seen that this result follows an analysis of SPAQ [**?**], which demonstrates that sharpness is the most highly related attribute to image quality.

Moreover, as we denoted at Sec 4.5, we conduct an additional experiment to verify that each attribute head accurately captures its corresponding attributes. To provide a detailed explanation, we manipulate each attribute by sequentially attaching relevant distortions, described in Sec B, to a given image and assess the significance of features by measuring the Grad-CAM from the MOS prediction to each feature map. Furthermore, to preserve the integrity of each feature map, we conduct this experiment using a linear probing setup. As shown in Figure 5, we first analyze the tendency of MOS variation in response to changes in image attribute. For sharpness and noisiness, we can observe that MOS decreases as more noise and blur are applied to the image. For contrast and brightness, the highest MOS is achieved when the attributes are at an optimal medium level, with the MOS decreasing when the attributes become either too high or too low. A similar result is observed for colorfulness, but the pattern of overall variation exhibits a slight difference. Interestingly, a negative correlation between the Mean Opinion Score (MOS) and Grad-CAM is observed for every attribute. It suggests that each attribute has a more significant impact on MOS prediction when a particular attribute varies to the extent that compromises the image quality. This observation indicates that each attribute head of ATTIQA can effectively capture changes in specific attributes from the original image.
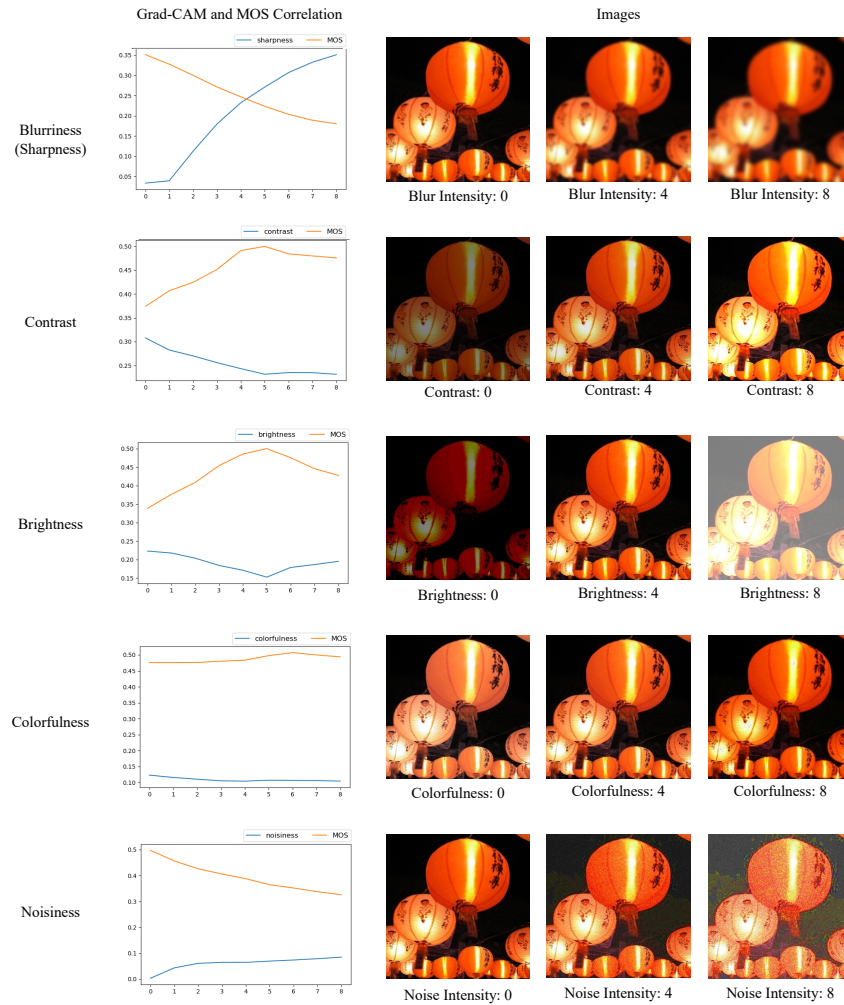
**Fig. 5:** Correlation graph illustrating the relationship between MOS and Grad-CAM across varying image attributes. On the right side, examples are provided to showcase how image attributes have been manipulated.

### C.3  Comparison with LIQE

In the main experiments, we only carried out single-dataset based experiments with LIQE to ensure a fair comparison with other baselines. To further assess the efficacy of our model in multiple-dataset environments, we measure ATTIQA's performances utilizing a dataset tailored for LIQE and its training recipe. As shown in Table 11, though ATTIQA is trained only using MOS, unlike LIQE which utilizes additional prompt-based annotations, it exhibits overall superior performances on given datasets. Exceptions are observed with two datasets, TID

**Table 11:** Performance comparison with LIQE.

| | Methods | Seen Dataset | | | | | | Unseen Dataset | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | LIVE | CSIQ | KADID | BID | LIVEC | KonIQ | TID | SPAQ | PIPAL |
| SR CC | LIQE | 0.970 | **0.936** | **0.930** | 0.875 | **0.904** | 0.919 | **0.811** | 0.881 | 0.478 |
| | ATTIQA | **0.974** | **0.936** | 0.923 | **0.891** | 0.901 | **0.920** | 0.796 | **0.891** | **0.521** |
| PL CC | LIQE | 0.951 | 0.939 | **0.931** | 0.900 | 0.910 | 0.908 | - | - | - |
| | ATTIQA | **0.977** | **0.950** | 0.927 | **0.918** | **0.917** | **0.932** | - | - | - |

and KADID, which primarily focus on specific distortions. These cases benefit from LIQE's approach of utilizing an additional dataset that contains additional annotations on synthetic distortions.

# D   Details on Application

## D.1   Metrics for Generative Model

This section provides details on the image generation pipeline used in a user study and more examples of results. For the generative model backbone, we utilize Stable Diffusion-2.1 [**?**], a widely used text-to-image synthesis model. To create diverse images, we establish prompt template as *"masterpiece, best quality, photorealistic photography, crystal clear, 8K UHD, {action} {object} {place}"* and generate images by replacing *"action"*, *"object"*, and *"place"* respectively with various candidates reported in Table 12. Figure 7 shows additional qualitative comparison results of a user study. Note that this user study aims to evaluate the quality of generated images solely by showing only the two images to the subject without any information on the corresponding prompts.

## D.2   Details on User Study

We conduct a user study using Amazon Technical Turk(AMT), gathering subjects in an anonymous setting without bias about gender or nationality. For the survey, the descriptions were presented as follows:

> 1) Participate in the image quality preference voting for tuning images. You can vote for the image you prefer between the two provided images
> 2) choose a better one with good image quality(color, sharpness, expressiveness..)

## D.3   Image Enhancement

In this section, we explain another application in ATTIQA's MOS that is used as a reward for reinforcement learning to find optimal parameters for the ISP [**?**]. The ISP pipeline consists of 14 modules, of which 16 tuning parameters of 7

**Table 12:** Prompt candidates utilized in image generation.

| Target | Prompt candidates |
|--------|-------------------|
| Action | *"playing", "running", "sleeping", "eating", "walking", "standing", "sitting", "jumping", "dancing"* |
| Object | *"a dog", "a cat", "a clothed man", "a dressed woman", "a bear"* |
| Place | *"in the grass", "in the room", "in the forest", "in the water", "in the snow", "in the desert"* |

modules were auto-tuned from the perspective of maximizing the ATTIQA's MOS. The tuned seven modules are as follows, and the numbers in parentheses of each module represent the number of tuning parameters: Tone-mapping (3), UV channel Denoising (3), Y channel Denoising (2), Brightness/Contrast control (2), Hue/Saturation control (2), Sharpening (3), and Texture enhancing (1). The agent of the reinforcement learning was trained using 730 raw images that were 4000x3000-sized. In the user survey, we uniformly sampled 200 raw images among 5000 raw images for fair comparisons. Figure 6 shows additional qualitative comparison results. It can also be seen that our pipeline retouches images to make them more colorful and vivid compared to both retouching by expert C and the default settings.
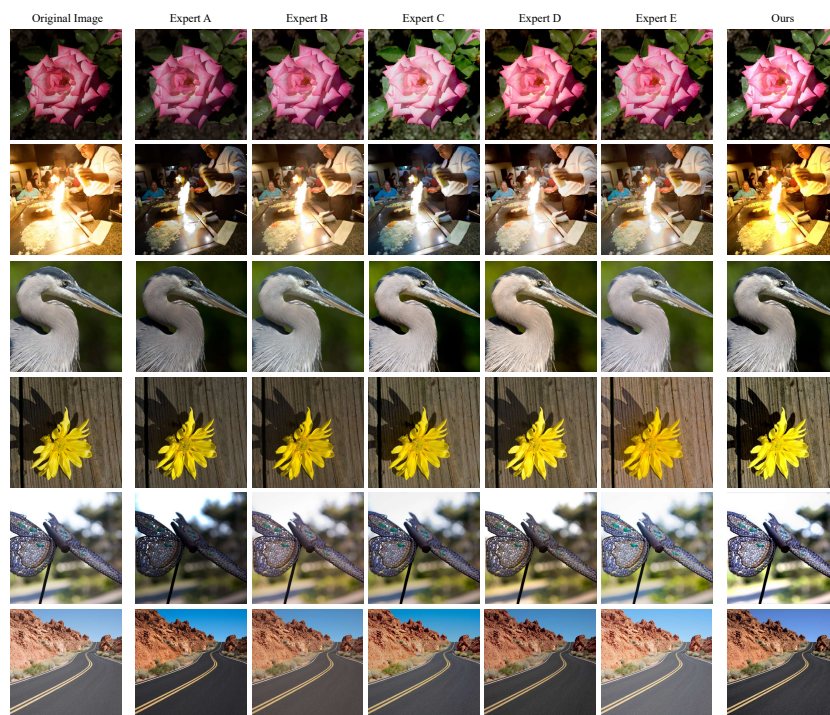
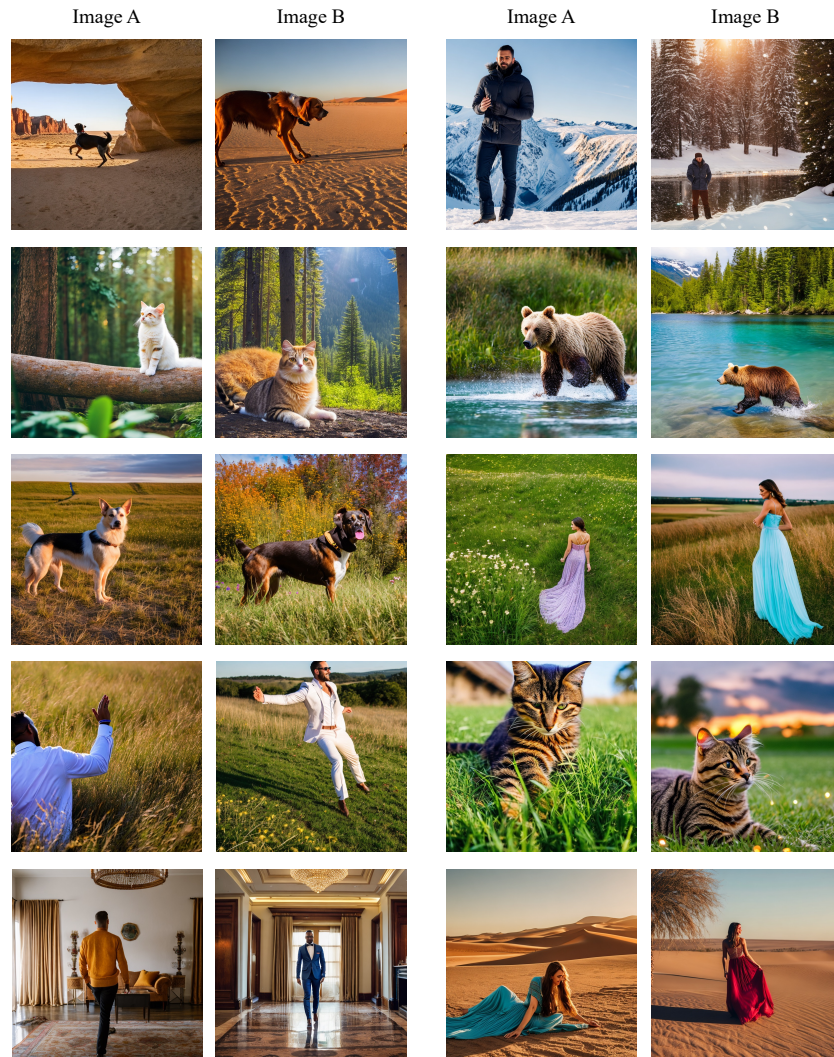**Fig. 6:** More qualitative comparisons about image enhancement.

**Fig. 7:** Examples of generated images. Each paired image is synthesized using the same prompt. Image A is the image that CONTRIQUE and Re-IQA predict the high quality score. In contrast, Image B is preferred by humans and assigned a high MOS by ATTIQA.

**Table 13:** Prompt candidates for each image attribute.

| Attribute | Prompt type | Prompt candidates |
|---|---|---|
| Sharpness | positive | "crisp","sharp","defined","clear","distinct","vivid","bright","detailed","refined","pristine", "flawless","lucid","exact","polished","pure","radiant","sleek","smooth","resolute","immaculate", "brilliant","vibrant","rich","clean","meticulous","unblemished","sublime","superior","splendid","exquisite", "true-to-life","tactile","textured","illuminated","lustrous","glossy","granular","pinpoint","spot-on","focused", "unambiguous","concise","intense","high-definition","lifelike","bold","harmonious","stunning","undistorted" |
| | negative | "blurry","fuzzy","hazy","vague","indistinct","muddled","obscured","smudged","cloudy","dull", "muted","out-of-focus","pixelated","jagged","noisy","grainy","mottled","muddy","murky","dim", "foggy","shadowy","bleary","washed-out","weak","gloomy","clouded","patchy","shrouded","veiled", "lacking clarity","rough","distorted","tarnished","ill-defined","ambiguous","flat","listless","pale","insipid", "smeared","streaked","stained","splotchy","spotty","blotchy","dingy","drab","tainted" |
| Contrast | positive | "crisp","defined","vivid","sharp","clear","distinguished","bold","pronounced","high-contrast","distinct", "lucid","striking","intense","robust","dynamic","stark","rich","deep","emphasized","highlighted", "vibrant","solid","notable","prominent","enhanced","powerful","contrastive","standout","bright","conspicuous", "discernible","marked","potent","compelling","dramatic","forceful","illuminated","brilliant","radiant","meticulous", "articulate","impressive","splendid","magnified","amplified","accentuated","divergent","outstanding","captivating" |
| | negative | "muddy","flat","blurred","faded","indistinct","washed-out","lackluster","dull","weak","subdued", "undistinguished","low-contrast","ambiguous","pale","muted","obscure","vague","insipid","clouded","nebulous", "tenuous","confusing","feeble","diminished","indiscernible","unnoticeable","slight","unemphasized","shadowed","doubtful", "hazy","unsaturated","ill-defined","unremarkable","bleak","insignificant","bland","monotone","uniform","muddled", "equivocal","unaccentuated","listless","understated","unimpressive","nondescript","faint","impotent","inaudible","discreet" |
| Brightness | positive | "luminous","bright","vivid","brilliant","radiant","gleaming","illuminated","clear","sparkling","shiny", "glowing","light-filled","dazzling","lucent","resplendent","shimmering","lustrous","beaming","crisp","vibrant", "intense","well-lit","brillante","glittering","glistening","blazing","effulgent","reflective","aglow","incandescent", "high-key","fiery","lambent","twinkling","opulent","sunlit","burnished","pristine","flashing","undimmed", "sunny","spotlit","blinding","flawless","translucent","glossy","crystal-clear","immaculate","gleamy" |
| | negative | "dim","dull","dark","shadowy","obscured","faint","gloomy","pale","muted","clouded", "bleak","underexposed","drab","faded","murky","shaded","veiled","flat","dusky","tenebrous", "sombre","gray","lackluster","washed-out","overcast","smoky","subdued","muffled","eclipsed","sullen", "unlit","opaque","low-key","blurred","darkened","blackened","shadowed","misty","lightless","moonless", "starless","inky","twilight","foggy","overclouded","cimmerian","umbrous","pitch-dark" |
| Colorfulness | positive | "vibrant","rich","vivid","saturated","brilliant","lively","radiant","bold","bright","colorful", "intense","resplendent","lush","deep","dazzling","varied","dynamic","electric","illuminated","vibrantly-hued","multicolored", "kaleidoscopic","strong","fluorescent","fiery","prismatic","stunning","flashy","beaming","pigmented","chromatic", "glistening","spectacular","polychromatic","sunny","iridescent","opulent","rainbow-like","effulgent","color-laden","invigorating", "gorgeous","lustrous","gleaming","dramatic","bursting","captivating","energetic" |
| | negative | "drab","dull","washed-out","muted","faded","pale","flat","monochrome","grayish","uninspiring", "lifeless","bleak","tarnished","insipid","blurred","cloudy","dim","neutral","colorless","lackluster", "subdued","murky","dusty","dusky","undistinguished","muddy","unsaturated","shadowed","overcast","veiled", "bland","indistinct","unvaried","uniform","faint","anaemic","vague","wan","stale","ashen", "pastel","watered-down","sallow","obscured","indeterminate","discolored","ill-defined","tinged","hazy" |
| Noisiness | positive | "clear","crisp","smooth","pure","sharp","pristine","flawless","noiseless","unblemished","intact", "clean","polished","immaculate","well-defined","impeccable","spotless","untainted","sleek","neat","unspoiled", "intact","refined","clutter-free","lucid","undisturbed","unmarred","untarnished","perfect","refreshing","distinct", "vivid","bright","detailed","accurate","faithful","exquisite","superb","top-notch","first-rate","matchless", "peerless","masterful","skillful","uncompromised","optimal","optimum","superior","prime","finest" |
| | negative | "grainy","speckled","mottled","patchy","dirty","blemished","marred","flecked","spotty","noisy", "blurry","fuzzy","hazy","cloudy","splotchy","streaky","dotted","gauzy","scratchy","scattered", "dappled","distorted","messy","smudged","lackluster","gloomy","dim","overcast","pockmarked","crackled", "choppy","erratic","spattered","discolored","inconsistent","irregular","shoddy","subpar","mediocre","unrefined", "vague","ambiguous","dull","dreary","stained","blemish-ridden","defective","imperfect","second-rate","tarnished","degraded" |