# Ambiguity of Objective Image Quality Metrics: A New Methodology for Performance Evaluation

Manri Cheon[a], Toinon Vigier[b], Lukáš Krasula[b], Junghyuk Lee[a], Patrick Le Callet[b], Jong-Seok Lee[a,*]

[a]*School of Integrated Technology, Yonsei University, 21983 Incheon, Korea*
[b]*LS2N UMR CNRS 6004, Université de Nantes, 44306 Nantes, France*

**Abstract**

Objective image quality metrics try to estimate the perceptual quality of the given image by considering the characteristics of the human visual system. However, it is possible that the metrics produce different quality scores even for two images that are perceptually indistinguishable by human viewers, which have not been considered in the existing studies related to objective quality assessment. In this paper, we address the issue of ambiguity of objective image quality assessment. We propose an approach to obtain an ambiguity interval of an objective metric, within which the quality score difference is not perceptually significant. In particular, we use the visual difference predictor, which can consider viewing conditions that are important for visual quality perception. In order to demonstrate the usefulness of the proposed approach, we conduct experiments with 33 state-of-the-art image quality metrics in the viewpoint of their accuracy and ambiguity for three image quality databases. The results show that the ambiguity intervals can be applied as an additional figure of merit when conventional performance measurement does not determine superiority between the metrics. The effect of the viewing distance on the ambiguity interval is also shown.

*Keywords:* Quality of experience, objective quality assessment, ambiguity interval, viewing distance

---

# 1. Introduction

Multimedia systems operating in resource-constrained environments usually strive to achieve two conflicting objectives: achieving efficiency and providing high quality content. For instance, compression, e.g., JPEG [2] and JPEG2000 [3] for images and H.264/AVC [4] and HEVC [5] for videos, is a representative way to deal with this issue; it can reduce the amount of data to enhance storage and transmission efficiency at the cost of degradation of perceptual quality. Quality degradation introduced through enhanced efficiency tends to lower the quality of experience (QoE) of the consumers. Therefore, it is important to carefully consider the trade-off relationship between the two objectives in designing the target multimedia systems and services.

The first step toward this goal is to accurately measure the perceptual quality of the content as perceived by human viewers, which can be performed via subjective quality assessment or objective quality assessment [6, 7, 8]. The former is the most accurate way of assessing the QoE, where human subjects are asked to rate the given content in terms of perceptual quality. However, it is time-consuming and expensive, and cannot be used in real time applications for controlling or optimizing the quality of the delivered content. Thus, objective quality assessment performed by objective metrics is widely used to replace subjective quality assessment, which tries to automatically predict perceived quality. A number of objective quality metrics have been developed and used for various applications including compression, transmission, enhancement, etc. [9].

It has been considered that the primary goal of an objective metric is to predict subjective quality scores, usually denoted as mean opinion scores (MOS), as accurately as possible. The ITU-T P.1401 standard [10] specifies recommended procedures to evaluate the accuracy of an objective quality metric. For instance, the Pearson's linear correlation coefficient (PLCC) and Spearman's rank ordered correlation coefficient (SROCC) are computed to evaluate linearity and monotonicity of metrics with respect to subjective data, respectively. In addition, the prediction error and consistency are also measured using the root-mean-square error (RMSE) and outlier ratio (OR), respectively. Additional statistical measures of performance have also been proposed in

2

[11].

In this paper, however, we argue that the accuracy is not the only perspective in which objective quality metrics should be judged, and propose that considering an additional figure of merit provides much more informative insight into the performance and behavior of the metrics, which is their *ambiguity* or, conversely, *reliability*. In general, the output of an objective metric for a given visual stimulus is expressed as a single value on a continuous scale. This means that when the predicted quality scores for a pair of stimuli by a metric are obtained, the quality superiority between the stimuli is always formed, no matter how small the difference is. However, a nonzero quality score difference between two similar stimuli may cause misleading conclusions when the quality difference is not perceivable by human viewers. In fact, the visual sensitivity of humans is limited in the sense that a small amount of pixel value difference is sometimes visually indistinguishable depending on several factors such as overall luminance and neighboring pixel values [12].

Figure 1 shows example images demonstrating the existence of ambiguity of objective metrics [13]. For two reference images (*parrots* and *house*) from the LIVE Image Quality Assessment Database [6], JPEG2000 compression is applied to corrupt them with different bitrates. When Figures 1(a) and 1(b) are visually compared, their quality difference can be easily perceived. We conducted a subjective quality assessment experiment using the paired comparison scheme [14, 15], where most of the hired subjects (14 out of 15) chose Figure 1(b) as the one having better quality. An objective metric, peak signal-to-noise ratio (PSNR), also rates Figure 1(b) as having better quality (with a difference of 2.49 dB), which is consistent with the quality superiority perceived by humans. On the other hand, the difference between Figures 1(c) and 1(d) is hardly noticeable; nearly a half of the subjects (6 out of 15) chose Figure 1(c). However, the quality measured by PSNR still determines that Figure 1(d) is better, showing a difference of 2.54 dB, which is even larger than the difference between Figures 1(a) and 1(b). Such inconsistent results between subjective and objective quality measurements are undesirable for quality-optimized multimedia systems. For instance, a system relying on PSNR may try to deliver Figure 1(d) instead of Figure 1(c) to improve QoE at the cost of an increased bitrate (20 to 35 kbytes), which is actually not so worthy for

3

(a) 30.46 dB        (b) 32.95 dB
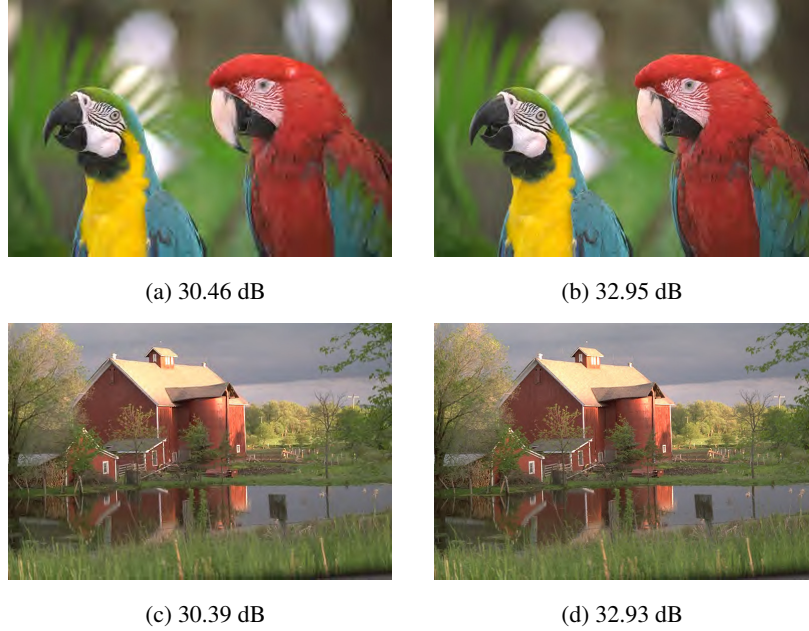
(c) 30.39 dB        (d) 32.93 dB

Figure 1:   Example images from the LIVE Image Quality Assessment Database [6], demonstrating the ambiguity of objective quality metrics (in this case, PSNR).

users. An additional observation in this example is the content-dependence of the ambiguity of objective metrics. In other words, the perceptual insignificance of the PSNR difference is observed only for *house*.

Even the state-of-the-art objective quality metrics showing good performance on predicting perceived quality (e.g., [16, 17, 18]) have the issue of indistinguishable quality ranges because all the existing metrics produce single numerical values representing the perceptual quality of given stimuli, which is highly related to the reliability of the metrics. In this paper, therefore, we address the issue of ambiguity of objective quality assessment and propose an approach to measure the ambiguity as an interval defining the indistinguishable quality score range, which can be applied to any quality metrics to supplement their usefulness in a new direction. Furthermore, we present use cases where the proposed approach can be useful, i.e., one for performance comparison of quality metrics' and the other for analysis of metrics performance in terms of reliability with respect to the viewing distance.

The main contributions of this paper can be summarized as follows:

1. **We propose an approach to measure the ambiguity of objective quality metrics.**

   The ambiguity is expressed as an interval on the scale of a metric's score, called *ambiguity interval*, within which the quality difference is perceptually indistinguishable. In obtaining ambiguity intervals, we incorporate the viewing conditions, in particular, viewing distance, because it is one of the most important factors that significantly influence the visual sensitivity of human viewers. Our approach employs the visual difference predictor (VDP) [19], which automatically estimates a threshold for perceptually indistinguishable pixel value difference at each pixel location. Using VDP also eliminates the necessity to conduct subjective experiments to obtain the ambiguity intervals, which maximizes the applicability of the proposed approach.

2. **We provide a practical use case, i.e., objective metric benchmarking, to demonstrate the effectiveness of the proposed approach.**

   We use the ambiguity characteristics of metrics for performance comparison of metrics in addition to the accuracy measure. It is shown that the ambiguity can play an important role to determine the superiority among the metrics. In the research community of multimedia quality assessment, systematic evaluation of objective metrics has been considered important to analyze their advantages and disadvantages [6, 7, 8, 20, 21, 22, 23]. The Video Quality Experts Group (VQEG), an international forum for perceptual quality assessment towards standardization, also puts a significant amount of efforts for this. Thus, this use case proposes a novel framework for benchmarking of objective quality metrics, which enables performance analysis of the metrics in multidimensional perspectives.

3. **As another practical use case, we evaluate state-of-the-art metrics in terms of viewing distance.**

   We show that the behavior of a metric depending on the viewing distance also provides valuable information in analyzing the metric's performance. Such information can be exploited as a part of benchmarking of objective metrics. In

addition, it can be used to identify proper viewing conditions where the metrics are reliable.

The rest of this paper is organized as follows. The following section presents the proposed approach in detail. Section 3 describes the experimental setup. The two use cases, where the ambiguity intervals are exploited, are given in Sections 4 and 5, respectively. Finally, conclusions are given in Section 6.

## 2. Proposed Method

### 2.1. Approach

As mentioned in the introduction, the goal of the proposed approach is to obtain an interval for a given objective quality metric, so that a score difference within the interval at that particular quality level is considered as being perceptually insignificant. The core idea to obtain such an interval is to change the amount of distortion (e.g., noise, compression artifacts, etc.) in an image and check using a perceptual model if the change of the distortion would be detected by human observers.

Algorithm 1 summarizes the procedure of the proposed approach to obtain the ambiguity interval (i.e., the upper and lower bounds of the interval) over the whole quality range for a source image and a type of distortion. Figure 2 illustrates the process to obtain the ambiguity interval for a particular quality level corresponding to a degraded image, which corresponds to lines 6 to 19 in Algorithm 1.

First, a quality degradation for the distortion type is applied to the source image ($I_0$) with various amounts of distortion and the objective quality levels of the resulting images are measured. Then, we determine perceptual distinguishability between two images having different amounts of artifacts. For a given image $I_i$ containing a certain type of artifacts, we obtain the level of ambiguity at the corresponding objective quality score ($Q_i$) as an interval around the score. We assess perceivable difference of the given image ($I_i$) compared to an image from the same source image but with different amounts of artifacts ($I_j$). We gradually increase (or decrease) the amounts of artifacts in $I_j$, until the images that are perceptually distinguishable from the given image are found. Among the images that are perceptually indistinguishable from $I_i$, the one with

6

**Algorithm 1** Computing the ambiguity interval

**Input:** Source image $I_0$ having $M$ pixels
**Output:** Upper bound width $U \in \mathbb{R}^N$ and lower bound width $L \in \mathbb{R}^N$ of the ambiguity interval

1: **for** $i \leftarrow 1, N$ **do**　　　　　　　　　　　$\triangleright$ $N$: number of considered quality levels
2:　　$I_i \leftarrow$ degrade_image$(I_0, i)$ $\triangleright$ Apply quality degradation (compression, blurring, etc.) to $I_0$ ($I_i$ is more degraded than $I_{i-1}$)
3:　　$Q_i \leftarrow$ measure_quality$(I_i)$　　　　$\triangleright$ Measure the objective quality (assume that a higher $Q_i$ indicates higher quality)
4: **end for**
5: **for** $i \leftarrow 1, N$ **do**
6:　　**for** $j \leftarrow i+1, N$ **do**
7:　　　　PMap $\leftarrow$ vdp$(I_i, I_j)$　　　　　　　　　$\triangleright$ Obtain the perceivableness map
8:　　　　**if** count(PMap $> 0.5)/M > k$ **then**
9:　　　　　　$L_i \leftarrow Q_i - Q_{j-1}$　　　　　　$\triangleright$ Obtain the width of the lower bound
10:　　　　　　break
11:　　　　**end if**
12:　　**end for**
13:　　**for** $j \leftarrow i-1, 1$ **do**
14:　　　　PMap $\leftarrow$ vdp$(I_i, I_j)$　　　　　　　　　$\triangleright$ Obtain the perceivableness map
15:　　　　**if** count(PMap $> 0.5)/M > k$ **then**
16:　　　　　　$U_i \leftarrow Q_{j-1} - Q_i$　　　　　　$\triangleright$ Obtain the width of the upper bound
17:　　　　　　break
18:　　　　**end if**
19:　　**end for**
20: **end for**
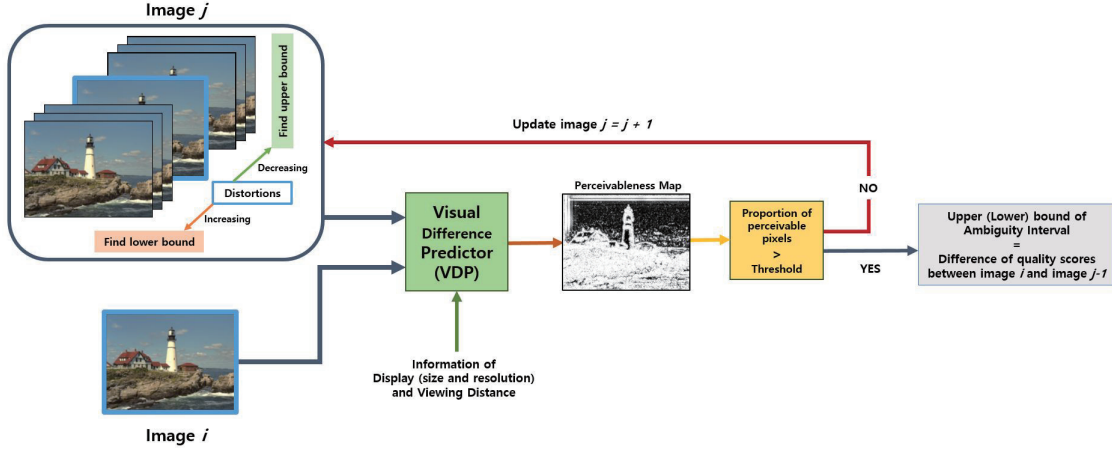21: **return** $U$ and $L$

Figure 2: Procedure to obtain an ambiguity interval based on a perceivableness map, which judges whether the two images are perceptually distinguishable or not. Note that the white pixels of the perceivableness map mean distinguishable pixels determined by VDP.
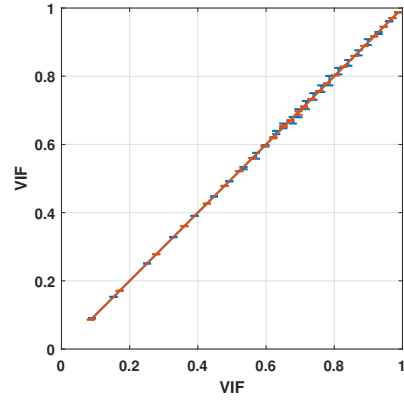
the highest (or lowest) quality level is identified, and the difference between the corresponding quality score and the quality score of $I_i$ is recorded as the width of the upper (or lower) bound of the interval, $U_i$ (or $L_i$).

A visual just-noticeable difference (JND) model is used to determine whether two images having different amounts of distortion are perceptually distinguishable. The JND model compares the two images and produces a map having the same size to that of the input images, called *perceivableness* map. Each pixel of the map represents the probability that the pixel value difference of the two images at the corresponding location is perceptually distinguishable. A probability of 0.5 (i.e., random chance) is considered as the threshold of distinguishability. Therefore, if at most a certain proportion (denoted as $k$) of the pixels of the perceivableness map have values above 0.5, the two images are considered to be perceptually indistinguishable.
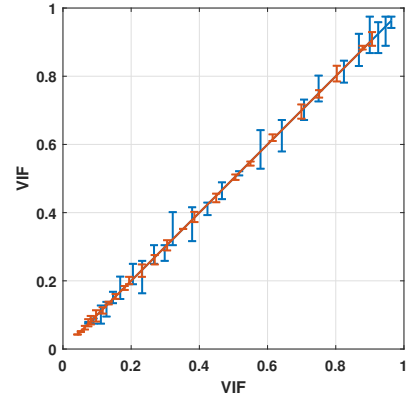
The JND model considered in this study is VDP, originally proposed by Daly [19]. It enables to specify the viewing conditions including the type, resolution, and parameters of the display, together with the viewing distance [24]. In particular, we use the

latest version, known as HDR-VDP 2.2 [16][1]. The model quantifies the visible difference between two input images under specific viewing conditions. The images are firstly passed through a model of the optical retinal pathway, including a simulation of intra-ocular light scatter, photoreceptor spectral sensitivity, luminance masking, and achromatic response. Further on, they are compared on multiple scales considering the model of neural noise, neural contrast sensitivity, and contrast masking. Note that when producing a perceivableness map, VDP takes into account the contextual information for each pixel (i.e., its relationship with neighboring pixels).
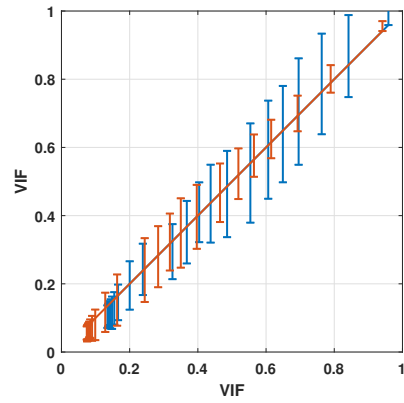
Figure 3 shows examples of the ambiguity intervals, which are obtained for the visual information fidelity (VIF) metric [17]. To determine the intervals, we generate $N$=100 images having different amounts of distortion (spanning the whole quality range) for each distortion type and each reference image in the LIVE Image Quality Assessment Database [6], and apply Algorithm 1 to them. In the figure, a higher score means a higher quality level, i.e., less artifacts. Three types of dependency of the interval are observed. First, the width of the interval is not necessarily uniform over the quality range. In Figure 3(c), for instance, the width of the interval is large for the intermediate quality range and small for low quality (near zero). This implies that the perceptual scale of the metric is not perfectly linear. Second, the interval width is dependent on the content, which is in line with the observation made from Figure 1. This is related to the fact that the visibility of quality degradation is dependent on the image content due to perceptual mechanisms such as frequency-dependent contrast sensitivity, spatial masking, etc. Third, the type of distortion also influences the interval because the detectability of quality difference depends on the type of artifacts. Detailed analysis is given in Section 4. In summary, the interval is dependent on the visual components included in the image, which are affected differently by the quality level, the distortion type, and the content itself.
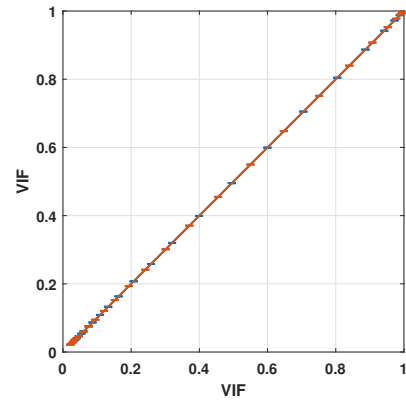
9

Figure 3: Examples of obtained ambiguity intervals of VIF for the LIVE database. The upper and lower bounds for two different reference images are expressed in different colors. (a) JPEG (b) JPEG2000 (JPEG2K) (c) Gaussian blur (GB) (d) white Gaussian noise (WN)
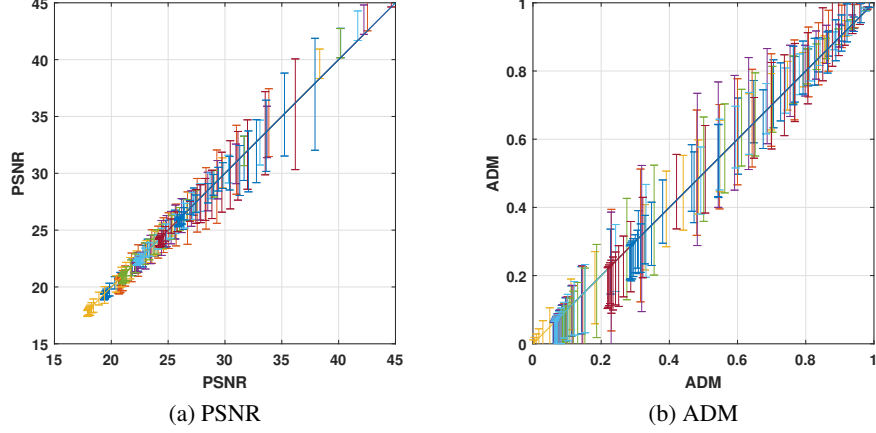
Figure 4: Examples of obtained ambiguity intervals for GB of the LIVE database. Different colors mean different reference images. (a) PSNR (b) ADM

## 2.2. Measures for Ambiguity Intervals

The ambiguity intervals of an objective metric can be used to measure the performance of the metric in terms of quality resolution. Figure 4 shows examples for two different metrics, i.e., PSNR and additive impairment and detail loss measure (ADM) [25], which have different output ranges and ambiguity interval widths. Overall, for instance, the intervals of ADM are larger than those of PSNR; the intervals of PSNR are relatively small for the low quality range and get larger as the quality increases, whereas the intervals of ADM are more uniform over the whole range. To enable easy comparison between the intervals of different metrics, we compute measures that summarize the ambiguity intervals of a metric. As the first step, the ambiguity intervals of a metric are normalized with the obtained output range of the metric, since different metrics may have different ranges and units. Note that in our preliminary work [1], nonlinear regression using subjective rating data was employed for normalization, which limits the applicability of the method only to the cases where subjective data are available. In addition, only the quality levels associated with subjective ratings were used, which permitted ambiguity evaluation only at a coarse level.

---

[1]An implementation is publicly available at http://hdrvdp.sourceforge.net/wiki/

11

Table 1: Characteristics of the three databases used for the experiments. A viewing distance is expressed as a multiple of the height of the display.

| Database | # Contents | Image resolutions | Distortion types | Screen | Viewing distance | Subjective ratings |
|---|---|---|---|---|---|---|
| LIVE [6] | 29 | 768×512 512×512 | JPEG, JPEG2K GB, WN | sRGB, CRT 21-inch, 1024×768 | 2H | DMOS |
| VDID [26] | 8 | 768×512 512×512 | JPEG, JPEG2K GB, WN | sRGB, LCD 23-inch, 1920×1080 | 4H 6H | DMOS |
| CIDIQ [27] | 23 | 800×800 | JPEG, JPEG2K GB, PN | sRGB, LCD 24-inch, 1920×1080 | 1.5H 3H | MOS |

We propose to compute three statistics of the ambiguity intervals, namely, the mean, maximum, and standard deviation of the widths of the ambiguity intervals over the whole quality range in order to measure the performance of a metric in multiple aspects of ambiguity. They are measures of the sensitivity of a metric in an average sense, the coarsest quality resolution, and the uniformity of the quality resolution, respectively. The smaller each of these measures is, the better the performance of the metric is.

## 3. Experimental setup

We conduct experiments in order to demonstrate applications where the proposed approach can be exploited effectively, which are shown in the following sections. This section explains the employed databases and the objective metrics considered in the experiments.

### 3.1. Databases

We employ three databases that are popularly used in the research of perceptual quality assessment, i.e., the LIVE Image Quality Assessment Database (LIVE) [6], which is one of the most popular databases for benchmarking objective metrics, the Viewing Distance-changed Image Database (VDID) [26], which is the first image quality assessment database specifically established for varying viewing distances, and the Colourlab Image Database: Image Quality (CIDIQ) [27], which also contains subjective data for multiple viewing distances. The databases were produced based on different experimental setups such as reference images, distortion types, screens, viewing distances, etc. We select them to ensure reproducibility of distortion types and availability of information regarding viewing environments, e.g., information of the screen

12

and viewing distance. Table 1 summarizes the characteristics of the databases. Four common distortion types are selected, i.e., JPEG compression, JPEG2000 (JPEG2K) compression, Gaussian blur (GB), and white Gaussian noise (WN). For the CIDIQ database, Poisson noise (PN) is considered instead of WN. JPEG and JPEG2K are well known compression schemes for images, and GB and WN (or PN) are distortions that can easily occur in pre- or post-processing of images. VDID and CIDIQ have subjective results from two different viewing distances.

### 3.2. Objective metrics

We consider 33 state-of-the-art objective quality metrics (28 full-reference (FR) metrics, one reduced-reference (RR) metric, and four no-reference (NR) metrics) for benchmarking. The tested FR metrics are PSNR, structural similarity index (SSIM) [18], multi-scale structural similarity (MS-SSIM) [28], visual signal-to-noise ratio (VSNR) [29], VIF [17], universal image quality index (UQI) [30], information fidelity criterion (IFC) [31], noise quality measure (NQM) [32], weighted signal to noise ratio (WSNR) [32], modified versions of PSNR (PSNR-HVS [33], PSNR-HVS-M [34], PSNR-HMA, PSNR-HA, PSNR-HMA-C, and PSNR-HA-C [35]), optimal scale selection (OSS)-PSNR and OSS-SSIM [26], information content weighted SSIM (IW-SSIM) [36], feature similarity index (FSIM) and chrominance extension of FSIM (FSIM-C) [37], gradient magnitude similarity deviation (GMSD) [38], most apparent distortion (MAD) [39], ADM [25], analysis of distortion distribution-based (ADD)-SSIM [40], ADD-gradient similarity index (ADD-GSIM) [40], a visual saliency-induced index (VSI) [41], image quality assessment based on gradient similarity (GSM) [42], and perceptual similarity (PSIM) [43]. The RR metric is reduced reference entropic differencing index (RRED) [44], and the NR metrics are spatial-spectral entropy-based quality (SSEQ) [45], oriented gradients image quality assessment (OG-IQA) [46], blind image integrity notator using DCT statistics (BLIINDS2) [47], and accelerated screen image quality evaluator (ASIQE) [48].

13

## 4. Use case 1 : Benchmarking of objective metrics

Objective quality metrics that can automatically predict perceived quality of visual content are a key component of quality-optimized multimedia systems. For instance, a method enhancing a given degraded image requires an objective metric as a criterion with respect to which the image is enhanced. Therefore, it is critical to identify a quality metric that mimics the human visual system as closely as possible, so that the results of optimization based on the metric are also optimal for human viewers. In this context, benchmarking studies of objective quality metrics have been conducted extensively in literature, e.g., [8, 22, 49, 50]. In these studies, as mentioned in the introduction, the prediction accuracy of existing metrics is considered as the most important performance index, which is typically measured in terms of PLCC, SROCC, OR, and RMSE. However, different metrics have different levels of ambiguity, which can be captured by the proposed approach. The use case presented in this section demonstrates how such information can be effectively used in the benchmarking.

In this use case, we use the LIVE database. The accuracy performance of the 33 state-of-the-art objective metrics is measured by PLCC between the ground truth subjective quality scores and the predicted quality scores[2]. In particular, PLCC is computed after nonlinear regression using the monotonic logistic function:

$$Q' = \beta_1 + \frac{\beta_2 - \beta_1}{1 + e^{-\left(\frac{Q - \beta_3}{\beta_4}\right)}} \tag{1}$$

to fit the objective scores outputted by a metric to the subjective quality scores, as described in the recommendation [51]. Here, $Q$ and $Q'$ denote the objective scores before and after regression, respectively. The initial values of the parameters ($\beta_1$ to $\beta_4$) are set as suggested in [51]. In addition, the statistical tests are also conducted [10], i.e., Z-tests are performed using the Fisher z-transformation for PLCC. The ambiguity performance of the metrics is evaluated based on the proposed approach. The mean, maximum, and standard deviation of the widths of the ambiguity intervals are

---

[2]Other measures such as SROCC, OR, and RMSE can be also used, but we use only PLCC for conciseness of presentation.
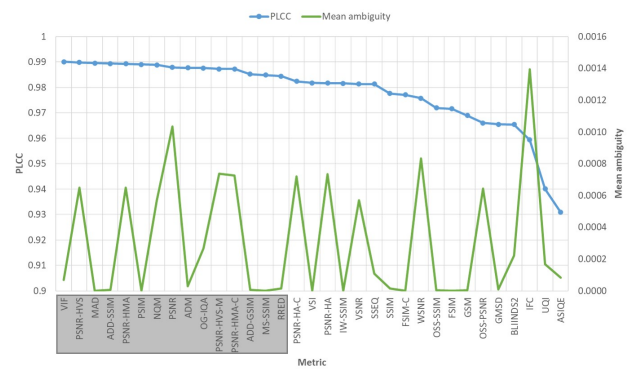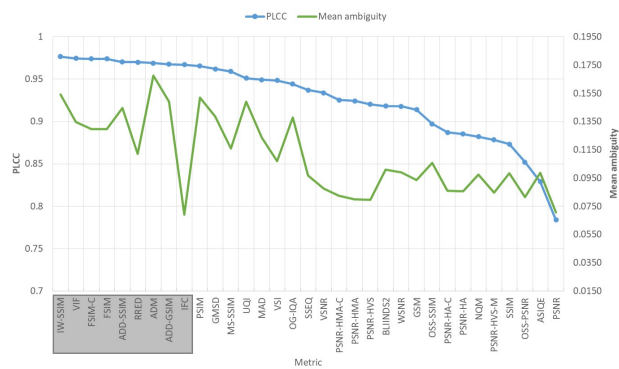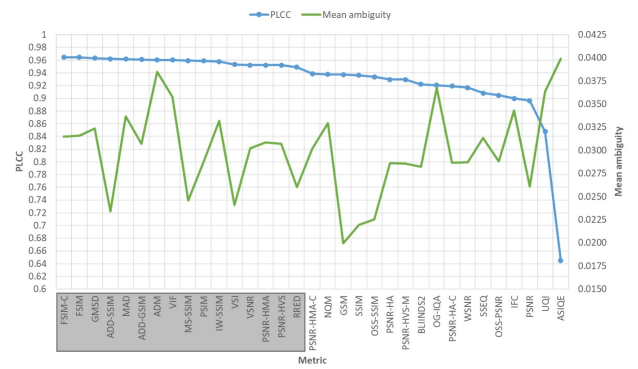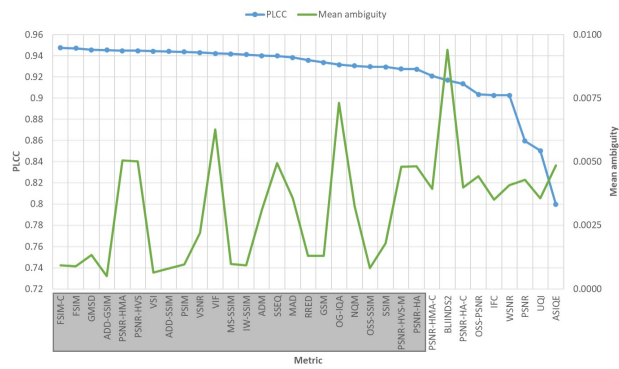
(a) JPEG

(b) JPEG2K

(c) GB

(d) WN

Figure 5: Performance of the objective metrics in terms of Pearson's linear correlation coefficient (PLCC) scores (blue) and mean of ambiguity intervals (green) for the LIVE database. (a) JPEG, (b) JPEG2K, (c) GB, and (d) WN. The metrics are listed in a descending order of the PLCC scores. The statistically equivalent metrics with the best metric for PLCC are marked in a gray box.

obtained. In addition, non-parametric Wilcoxon-Mann-Whitney tests are conducted to statistically compare the ambiguity intervals of different metrics.

Figure 5 summarizes the PLCC values and the mean ambiguity intervals of the 33 metrics. The results for the four distortion types are shown separately, and the metrics are listed in a descending order of the PLCC values. In the figure, the metrics showing statistically equivalent performance with the best metric in terms of PLCC are marked in the gray box. We can observe that the superiority of a metric over the others in terms of accuracy may not coincide with its superiority in terms of ambiguity, and vice versa. For instance, in Figure 5(b), the best metric in terms of accuracy is FSIM-C, but GSM,

15

which is statistically significantly inferior to FSIM-C, is the best in terms of ambiguity.

Many metrics predict perceived image quality with high accuracy. For instance, the best metric in terms of PLCC for JPEG in Figure 5(a), i.e., FSIM-C, which shows PLCC of about 0.95, is not statistically different with PSNR-HA, which ranks 24th. Thus, it would be difficult to distinguish the superiority between these metrics. At this point, we can apply the results of the ambiguity analysis. Among the top 24 metrics, ADD-GSIM has the smallest mean width of the ambiguity intervals, which is revealed to be significantly smaller than the second smallest one (VSI) by the statistical test ($p < 0.01$). For the other types of distortion, similar trends are also observed, i.e., a number of the metrics show similar performance in terms of PLCC and their performance is not statistically different from that of the best metric, and we can use ambiguity intervals in order to choose the best metric for these cases. From this approach, ADD-SSIM, IFC, and MAD are selected as the best metrics for JPEG2K, GB, and WN, respectively.

The mean ambiguity interval widths are different depending on the distortion type. The average values for all metrics are 0.0032, 0.0300, 0.1104, and 0.0003 for JPEG, JPEG2K, GB, and WN, respectively. The smallest ambiguity intervals are produced for WN, because changes of the amount of white noise can be easily detected compared to the other types of distortion. The GB distortion yields the largest ambiguity intervals because the change of the strength of GB is relatively hard to distinguish. JPEG2K also has relatively large ambiguity intervals because the introduced artifacts in the images are quite similar to those by GB.

In addition to the mean width of the ambiguity intervals, the maximum and standard deviation of the intervals can be also considered in order to analyze the performance of metrics and find the superiority between them. Figure 6 shows the performance of the objective metrics for all distortion types, which have statistically equivalent performance in terms of PLCC. When we compare the metrics, e.g., SSEQ and IW-SSIM, the two metrics show similar performance based on PLCC and the mean ambiguity interval. The maximum and standard deviation of the ambiguity intervals are smaller for IW-SSIM, which can be regarded as a better metric; a low standard deviation of the ambiguity intervals means that it has a uniform quality resolution (or ambiguity) for all quality ranges, which is useful in applications where the metric needs to operate
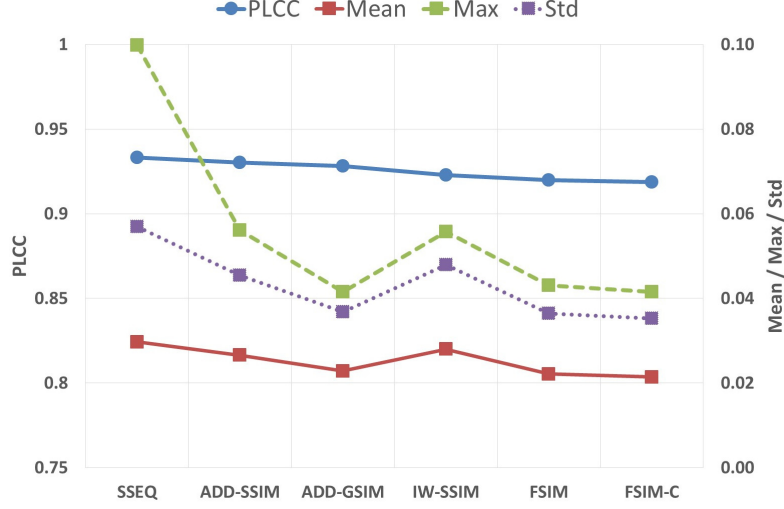
Figure 6: Performance of the top-performing objective metrics showing statistically equivalent PLCC values for all data of the LIVE database. PLCC scores and the mean, maximum, and standard deviation values of ambiguity intervals are shown.

in a wide range of quality. As another example, ADD-SSIM, ADD-GSIM, FSIM, and FSIM-C show statistically equivalent performance in terms of the mean ambiguity intervals, showing mean ambiguity intervals of only about 2.0-2.5% of the whole quality range. However, the maximum and standard deviation of the ambiguity intervals of ADD-SSIM are larger than those of the other three metrics, and thus it may be less preferable. Therefore, considering all the ambiguity measures, ADD-GSIM, FSIM, and FSIM-C can be regarded as the best metrics.

## 5. Use case 2 : Viewing distance vs. ambiguity

The viewing distance is one of the most important factors that influence visual quality perception of human viewers. As the distance from a viewer to an image gets large, less and less details in the image are distinguished, changes or artifacts in the image become less noticeable, and the viewer's quality perception becomes less reliable. The proposed approach incorporates this tendency by employing the VDP method that considers the viewing environment including the viewing distance.

However, due to the difference in underlying mechanisms of different objective
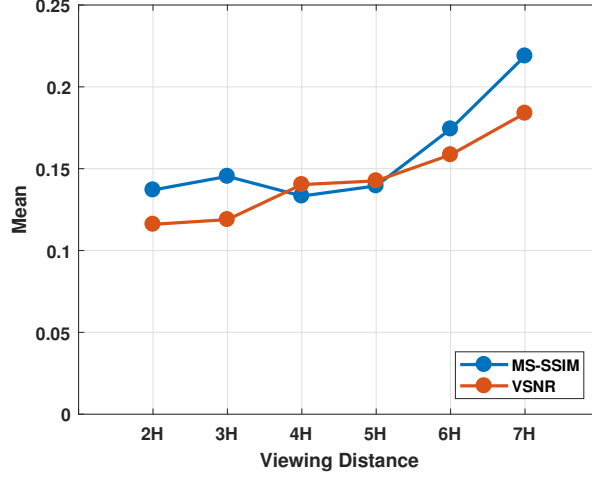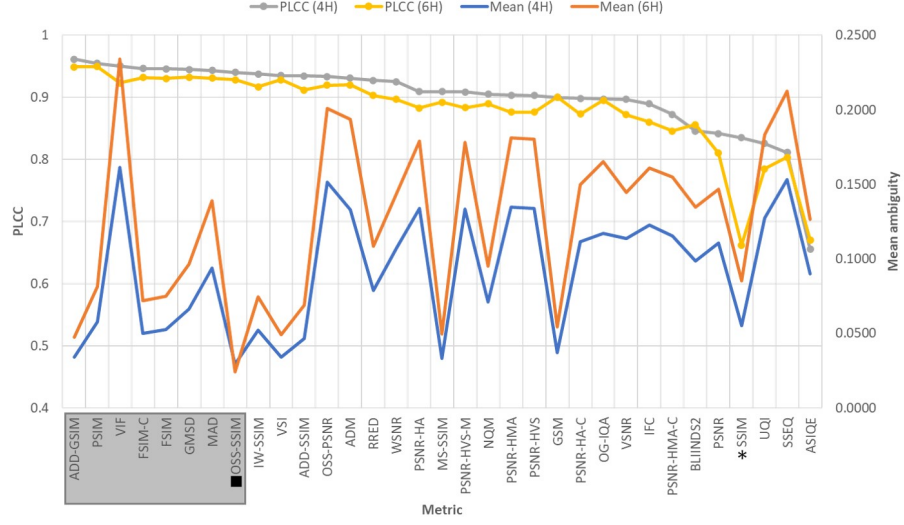
17

Figure 7: Examples of ambiguity intervals of two objective metrics, MS-SSIM (blue) and VSNR (orange), with respect to the viewing distance (in multiples of the display height) for GB of the VDID database. The superiority of a metric against the other varies depending on the viewing distance.

metrics, they may show different ambiguity patterns with respect to the viewing distance. For instance, the superiority of the metrics in terms of the ambiguity interval may change depending on the viewing distance. Figure 7 shows the mean ambiguity intervals of two metrics, MS-SSIM and VSNR, for GB of the VDID database with respect to the viewing distance. When the viewing distance is 4 or 5 times the display height (i.e., 4H or 5H), MS-SSIM shows slightly smaller ambiguity intervals than VSNR, whereas VSNR shows smaller intervals than MS-SSIM for the other viewing distances. Thus, the viewing distance should be considered carefully when the ambiguity of a metric is evaluated. In general, it is preferable for a metric not only to have high accuracy and low ambiguity for a particular viewing distance, but also to show consistent performance over various viewing distances in terms of both accuracy and ambiguity.
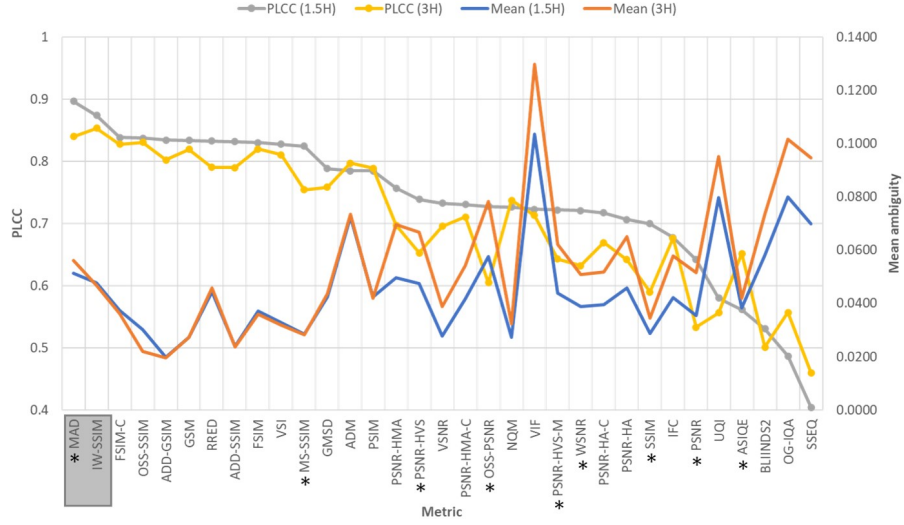
In this section, we demonstrate that the ambiguity behavior of metrics with respect to the viewing distance can be used to compare the reliability performance of the metrics, which can be seen as an extension of the benchmarking in the previous section, and to identify proper viewing distances for which a metric can be used reliably. The VDID and CIDIQ databases are used.

Performance of the metrics for two viewing distances in terms of PLCC and mean of the ambiguity intervals is shown in Figure 8. Most of the metrics show statistically equivalent PLCC scores for the two viewing distances; only one and nine metrics show significantly different accuracy scores for VDID and CIDIQ, respectively (which are marked with asterisks in Figure 8). However, for VDID, all metrics except for OSS-SSIM (marked with a square in Figure 8(a)) show significantly different ambiguity interval widths for the two viewing distances. Furthermore, OSS-SSIM shows high accuracy, i.e., it is included in the group of top-performing metrics (showing statistically equivalent PLCC scores with the best one for the short distance), and shows the smallest mean ambiguity intervals for both viewing distances (which are statistically equivalent). Thus, we can choose OSS-SSIM as the best metric considering both the accuracy and the ambiguity for different viewing distances. OSS-SSIM explicitly considers the effect of the viewing distance, which seems to be the reason for the consistency of its ambiguity performance. In the case of CIDIQ, all metrics have significantly different results of ambiguity intervals. MAD and IW-SSIM are two top-performing metrics in terms of accuracy for the short distance. However, these metrics have relatively lower performance in terms of ambiguity (i.e., larger mean interval widths) than the following ones (in the ranking of accuracy), e.g., OSS-SSIM and ADD-GSIM. If we accept a slight loss in terms of accuracy, it would be a better choice to select ADD-GSIM or OSS-SSIM as the best metric with consideration of both the accuracy and ambiguity for the two viewing distances.

Next, we analyze patterns of the ambiguity intervals over various viewing distances. As an example, Figure 9 shows the mean widths of the ambiguity intervals of ADD-GSIM for each of the four distortion types of VDID. As aforementioned, as the viewing distance increases, the ability of human viewers to distinguish the details in images decreases. The ambiguity intervals obtained by our approach also tend to increase with the increasing viewing distance. A gradual increase of the ambiguity intervals due to increase of the viewing distance is acceptable, but a sudden increase of the slope would not be desirable. For instance, in Figure 9(c), the slope for GB increases suddenly after 5H, thus, care must be taken when the metric is used for viewing distances larger than 5H.

(a) VDID



(b) CIDIQ

Figure 8: Performance of the metrics for two viewing distances in terms of PLCC scores and mean of ambiguity intervals for the (a) VDID and (b) CIDIQ databases. The metrics are listed in a descending order of the PLCC for the short viewing distances (4H for VDID and 1.5H for CIDIQ). The metrics having statistically different accuracy between the two viewing distances are marked with asterisks. The statistically equivalent metrics with the best metric in terms of PLCC for the short viewing distance are marked with a gray box. The metric having statistically equivalent ambiguity interval widths for the two viewing distances is marked with a square.

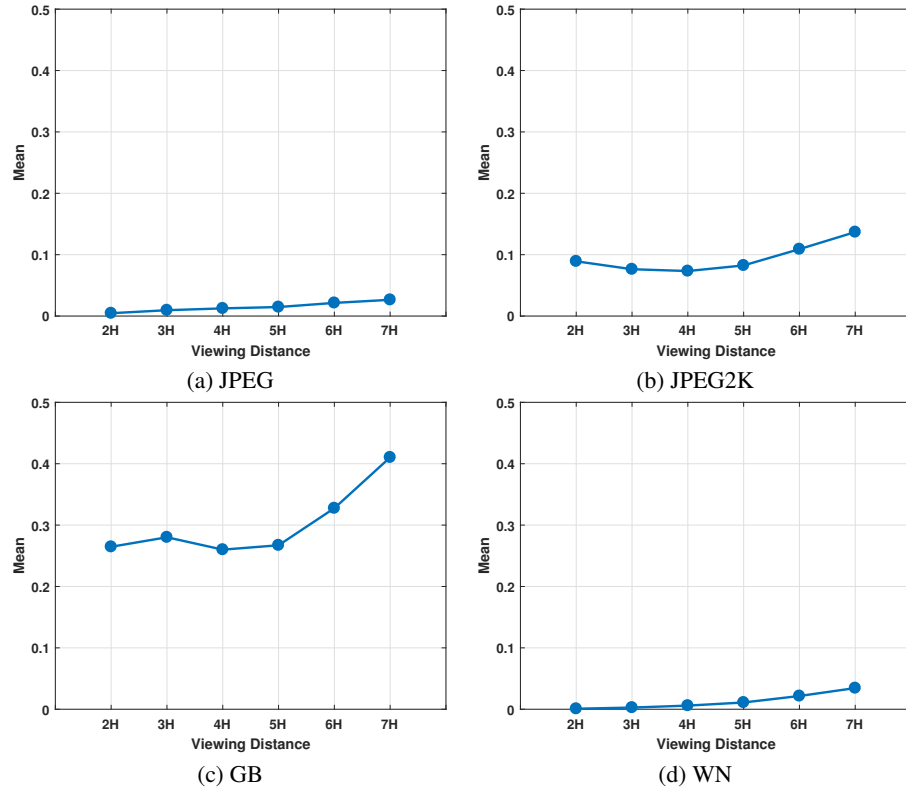(a) JPEG

(b) JPEG2K

(c) GB

(d) WN

Figure 9: Mean widths of the ambiguity intervals of ADD-GSIM for the VDID database (a) JPEG, (b) JPEG2K, (c) GB, and (d) WN. The distortion type influences the slopes of the curves.

Figures 10 and 11 show the mean ambiguity interval widths of the metrics with respect to the viewing distance for VDID and CIDIQ, respectively. For each distortion type, the metrics are sorted with respect to the mean ambiguity interval width for the short viewing distances, and the result of a metric in each quarter is presented. We can observe that the overall slopes of the mean ambiguity interval due to the viewing distance change are different depending on the objective metrics, distortion types, and databases.

When we compare the metrics for the same distortion types (i.e., the four panels in each row), the worse the performance of the metric in terms of ambiguity is, the larger the slope of the graph is. This tendency is observed clearly except for GB of CIDIQ (Figure 11(c)). Therefore, choosing a metric showing good performance in terms of ambiguity for a particular viewing distance is useful also for its reliable usage over different viewing distances.

In Figures 10 and 11, it is also observed that the shapes of the graphs are different depending on the distortion type. The cases of JPEG and JPEG2K for both databases show mostly monotonically increasing patterns. Monotonic increases are also observed for the two types of noise (WN of VDID and PN of CIDIQ), except for the clipping at zero for PN of CIDIQ. The graphs of GB show different tendencies; the mean widths of the intervals remain almost the same for some ranges of viewing distance (for small distances up to 5H for VDID and all distances considered for CIDIQ). As mentioned earlier, when the viewing distance increases, an ability to distinguish the details in the image also decreases. Since GB has already reduced the details in the image, the ambiguity intervals are less affected by the viewing distance change. In some cases, there exist sudden increases of the ambiguity interval widths (GB of VDID and PN of CIDIQ, both after 5H), which needs to be carefully considered when a system using a quality metric operates for a wide range of viewing distances.

## 6. Conclusion

In this paper, we have proposed a new way to measure performance of objective image quality metrics in the viewpoint of the quality resolution. The procedure to ob-
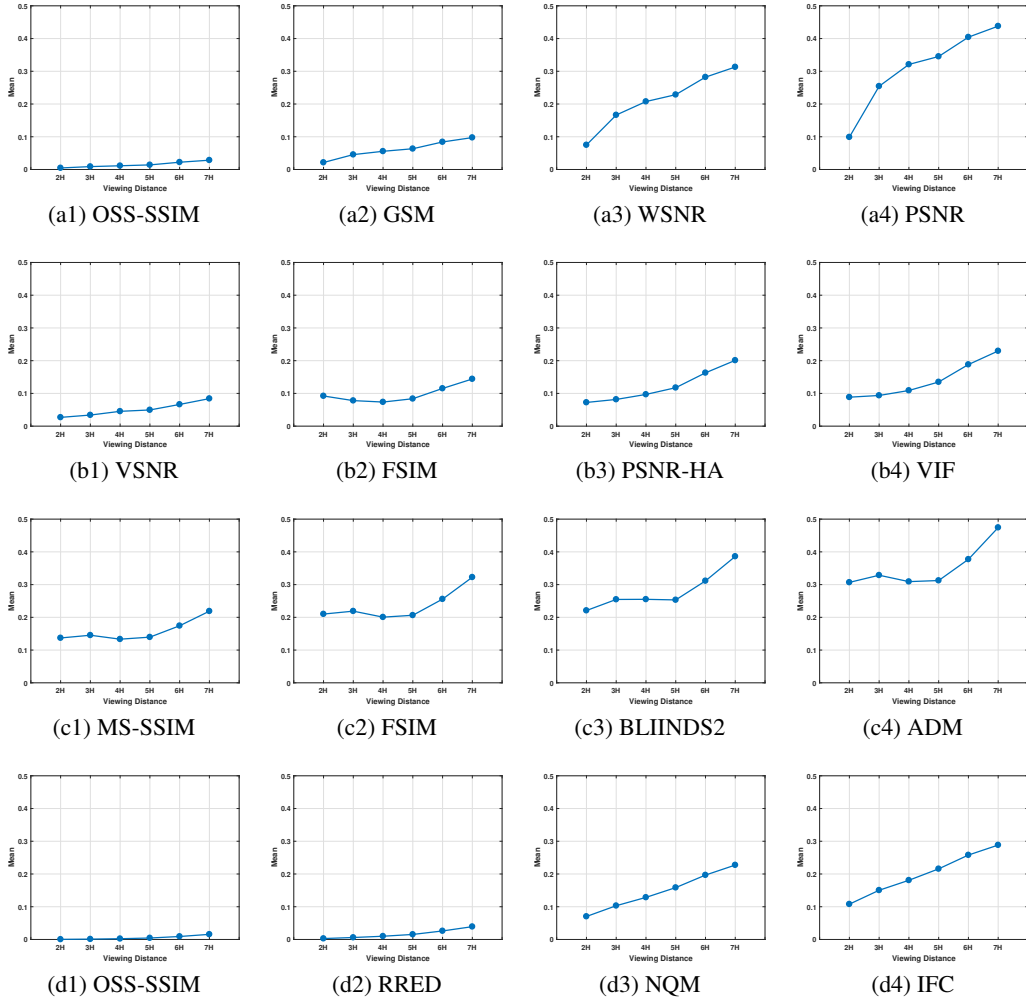
Figure 10: Mean widths of the ambiguity intervals of the objective metrics for the VDID database. For each distortion type, metrics in the first, second, third, and last quarters in the ascending order of the mean ambiguity interval width for 4H are shown from left to right. (a1)-(a4) JPEG, (b1)-(b4) JPEG2K, (c1)-(c4) GB, and (d1)-(d4) WN
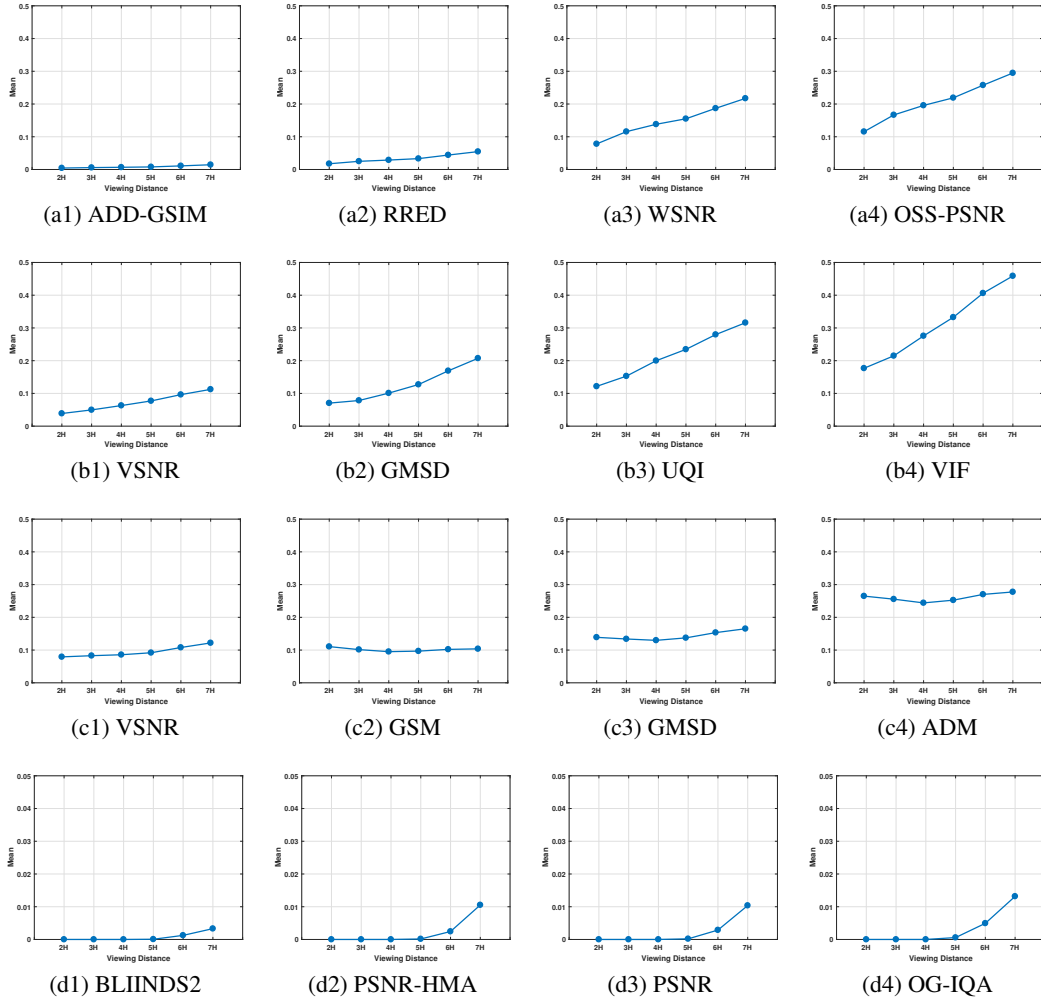
Figure 11: Mean widths of the ambiguity intervals of the objective metrics for the CIDIQ database. For each distortion type, metrics in the first, second, third, and last quarters in the ascending order of the mean ambiguity interval width for 1.5H are shown from left to right. (a1)-(a4) JPEG, (b1)-(b4) JPEG2K, (c1)-(c4) GB, and (d1)-(d4) PN

tain the ambiguity interval for an objective quality score has been developed. We have demonstrated that the width and the uniformity of the interval over the quality range are useful as performance measures in addition to the accuracy of quality estimation for comparison of different metrics. In addition, the need for consideration of the viewing distance has been emphasized when the ambiguity intervals are used.

In addition to the use cases shown in this paper, there are several other potential applications of the proposed method. An example is to construct rate-distortion (R-D) curves having ambiguity intervals for evaluating image compression methods, where the distortion is measured with an objective metric and the ambiguity interval at each rate value is obtained using the proposed method, which is an objective counterpart of the method obtaining R-D curves having intervals based on subjective quality scores [52].

One possible follow-up research question is: Is there a way to combine the two (possibly conflicting) performance dimensions of objective quality metrics (i.e., accuracy and ambiguity) to have a single performance measure? A general solution to this issue may be very challenging to develop. However, some guidelines to define the superiority and inferiority among different metrics could be identified depending on the target application, which would be desirable to explore as future work.

**Acknowledgment**

**References**

[1] M. Cheon, J.-S. Lee, Ambiguity-based evaluation of objective quality metrics for image compression, in: Proc. Int. Conf. Quality of Multimedia Experience

(QoMEX), 2016, pp. 1–6.

[2] G. K. Wallace, The JPEG still picture compression standard, Communications of the ACM 34 (4) (1991) 30–44.

[3] A. Skodras, C. Christopoulos, T. Ebrahimi, The JPEG 2000 still image compression standard, IEEE Signal Processing Magazine 18 (5) (2001) 36–58.

[4] T. Wiegand, G. J. Sullivan, G. Bjontegaard, A. Luthra, Overview of the H.264/AVC video coding standard, IEEE Trans. Circuits and Systems for Video Technology 13 (7) (2003) 560–576.

[5] G. J. Sullivan, J. Ohm, W.-J. Han, T. Wiegand, Overview of the high efficiency video coding (HEVC) standard, IEEE Trans. Circuits and Systems for Video Technology 22 (12) (2012) 1649–1668.

[6] H. R. Sheikh, M. F. Sabir, A. C. Bovik, A statistical evaluation of recent full reference image quality assessment algorithms, IEEE Trans. Image Processing 15 (11) (2006) 3440–3451.

[7] S. Chikkerur, V. Sundaram, M. Reisslein, L. J. Karam, Objective video quality assessment methods: A classification, review, and performance comparison, IEEE Trans. Broadcasting 57 (2) (2011) 165–182.

[8] M. Cheon, J.-S. Lee, Subjective and objective quality assessment of compressed 4K UHD videos for immersive experience, IEEE Trans. Circuits and Systems for Video Technology 28 (7) (2018) 1467–1480.

[9] Z. Wang, Applications of objective image quality assessment methods [applications corner], IEEE Signal Processing Magazine 28 (6) (2011) 137–142.

[10] ITU-T, Recommendation P.1401: Methods, metrics and procedures for statistical evaluation, qualification and comparison of objective quality prediction models, Tech. rep., ITU-T (2012).

[11] L. Krasula, K. Fliegel, P. Le Callet, M. Klíma, On the accuracy of objective image and video quality models: New methodology for performance evaluation, in:

Proc. Int. Workshop on Quality of Multimedia Experience (QoMEX), 2016, pp. 1–6.

[12] N. Jayant, J. Johnston, R. Safranek, Signal compression based on models of human perception, Proc. IEEE 81 (10) (1993) 1385–1422.

[13] M. Cheon, J.-S. Lee, On ambiguity of objective image quality assessment, Electronics Letters 52 (1) (2016) 34–35.

[14] ITU-R, Recommendation BT.500-13: Methodology for the subjective assessment of the quality of television, Tech. rep., ITU-R (2012).

[15] J.-S. Lee, On designing paired comparison experiments for subjective multimedia quality assessment, IEEE Trans. Multimedia 16 (2) (2014) 564–571.

[16] M. Narwaria, R. K. Mantiuk, M. P. Da Silva, P. Le Callet, HDR-VDP-2.2: A calibrated method for objective quality prediction of high-dynamic range and standard images, Journal of Electronic Imaging 24 (1) (2015) 1–3.

[17] H. R. Sheikh, A. C. Bovik, Image information and visual quality, IEEE Trans. Image Processing 15 (2) (2006) 430–444.

[18] Z. Wang, A. C. Bovik, H. R. Sheikh, E. P. Simoncelli, Image quality assessment: from error visibility to structural similarity, IEEE Trans. Image Processing 13 (4) (2004) 600–612.

[19] S. J. Daly, Visible differences predictor: an algorithm for the assessment of image fidelity, in: Proc. SPIE/IS&T Symposium on Electronic Imaging: Science and Technology, 1992, pp. 2–15.

[20] P. Mohammadi, A. Ebrahimi-Moghadam, S. Shirani, Subjective and objective quality assessment of image: A survey, Majlesi Journal of Electrical Engineering 9 (1) (2015) 55–83.

[21] W. Lin, C.-C. J. Kuo, Perceptual visual quality metrics: A survey, Journal of Visual Communication and Image Representation 22 (4) (2011) 297–312.

[22] J.-S. Lee, Comparison of objective quality metrics on the scalable extension of H.264/AVC, in: Proc. IEEE Int. Conf. Image Processing (ICIP), 2012, pp. 693–696.

[23] M. Cheon, J.-S. Lee, Evaluation of objective quality metrics for multidimensional video scalability, Journal of Visual Communication and Image Representation 35 (2016) 132–145.

[24] R. Mantiuk, K. J. Kim, A. G. Rempel, W. Heidrich, HDR-VDP-2: a calibrated visual metric for visibility and quality predictions in all luminance conditions, ACM Trans. Graphics 30 (4) (2011) 40:1–13.

[25] S. Li, F. Zhang, L. Ma, K. N. Ngan, Image quality assessment by separately evaluating detail losses and additive impairments, IEEE Trans. Multimedia 13 (5) (2011) 935–949.

[26] K. Gu, M. Liu, G. Zhai, X. Yang, W. Zhang, Quality assessment considering viewing distance and image resolution, IEEE Trans. Broadcasting 61 (3) (2015) 520–531.

[27] X. Liu, M. Pedersen, J. Y. Hardeberg, CID:IQ–a new image quality database, in: Proc. Int. Conf. Image and Signal Processing, 2014, pp. 193–202.

[28] Z. Wang, E. P. Simoncelli, A. C. Bovik, Multiscale structural similarity for image quality assessment, in: Proc. IEEE Asilomar Conf. Signals, Systems and Computers, Vol. 2, 2003, pp. 1398–1402.

[29] D. M. Chandler, S. S. Hemami, VSNR: A wavelet-based visual signal-to-noise ratio for natural images, IEEE Trans. Image Processing 16 (9) (2007) 2284–2298.

[30] Z. Wang, A. C. Bovik, A universal image quality index, IEEE Signal Processing Letters 9 (3) (2002) 81–84.

[31] H. R. Sheikh, A. C. Bovik, G. De Veciana, An information fidelity criterion for image quality assessment using natural scene statistics, IEEE Trans. Image Processing 14 (12) (2005) 2117–2128.

[32] N. Damera-Venkata, T. D. Kite, W. S. Geisler, B. L. Evans, A. C. Bovik, Image quality assessment based on a degradation model, IEEE Trans. Image Processing 9 (4) (2000) 636–650.

[33] K. Egiazarian, J. Astola, N. Ponomarenko, V. Lukin, F. Battisti, M. Carli, New full-reference quality metrics based on HVS, in: Proc. Int. Workshop on Video Processing and Quality Metrics, Vol. 4, 2006.

[34] N. Ponomarenko, F. Silvestri, K. Egiazarian, M. Carli, J. Astola, V. Lukin, On between-coefficient contrast masking of DCT basis functions, in: Proc. Int. Workshop on Video Processing and Quality Metrics, Vol. 4, 2007.

[35] N. Ponomarenko, O. Ieremeiev, V. Lukin, K. Egiazarian, M. Carli, Modified image visual quality metrics for contrast change and mean shift accounting, in: Proc. Int. Conf. The Experience of Designing and Application of CAD Systems in Microelectronics, 2011, pp. 305–311.

[36] Z. Wang, Q. Li, Information content weighting for perceptual image quality assessment, IEEE Trans. Image Processing 20 (5) (2011) 1185–1198.

[37] L. Zhang, D. Zhang, X. Mou, FSIM: a feature similarity index for image quality assessment, IEEE Trans. Image Processing 20 (8) (2011) 2378–2386.

[38] W. Xue, L. Zhang, X. Mou, A. C. Bovik, Gradient magnitude similarity deviation: A highly efficient perceptual image quality index, IEEE Trans. Image Processing 23 (2) (2014) 684–695.

[39] E. C. Larson, D. M. Chandler, Most apparent distortion: full-reference image quality assessment and the role of strategy, Journal of Electronic Imaging 19 (1) (2010) 1–21.

[40] K. Gu, S. Wang, G. Zhai, W. Lin, X. Yang, W. Zhang, Analysis of distortion distribution for pooling in image quality prediction, IEEE Trans. Broadcasting 62 (2) (2016) 446–456.

[41] L. Zhang, Y. Shen, H. Li, VSI: A visual saliency-induced index for perceptual image quality assessment, IEEE Trans. Image Processing 23 (10) (2014) 4270–4281.

[42] A. Liu, W. Lin, M. Narwaria, Image quality assessment based on gradient similarity, IEEE Trans. Image Processing 21 (4) (2012) 1500–1512.

[43] K. Gu, L. Li, H. Lu, X. Min, W. Lin, A fast reliable image quality predictor by fusing micro- and macro-structures, IEEE Trans. on Industrial Electronics 64 (5) (2017) 3903–3912.

[44] R. Soundararajan, A. C. Bovik, RRED indices: Reduced reference entropic differencing for image quality assessment, IEEE Trans. Image Processing 21 (2) (2012) 517–526.

[45] L. Liu, B. Liu, H. Huang, A. C. Bovik, No-reference image quality assessment based on spatial and spectral entropies, Signal Processing: Image Communication 29 (8) (2014) 856–863.

[46] L. Liu, Y. Hua, Q. Zhao, H. Huang, A. C. Bovik, Blind image quality assessment by relative gradient statistics and adaboosting neural network, Signal Processing: Image Communication 40 (2016) 1–15.

[47] M. A. Saad, A. C. Bovik, C. Charrier, Blind image quality assessment: A natural scene statistics approach in the DCT domain, IEEE Trans. Image Processing 21 (8) (2012) 3339–3352.

[48] K. Gu, J. Zhou, J. Qiao, G. Zhai, W. Lin, A. C. Bovik, No-reference quality assessment of screen content pictures, IEEE Trans. on Image Processing 26 (8) (2017) 4005–4018.

[49] P. Hanhart, P. Korshunov, T. Ebrahimi, Benchmarking of quality metrics on ultra-high definition video sequences, in: Proc. Int. Conf. Digital Signal Processing (DSP), 2013, pp. 1–8.

[50] S. Tian, L. Zhang, L. Morin, O. Déforges, A benchmark of DIBR synthesized view quality assessment metrics on a new database for immersive media applications, IEEE Trans. Multimedia 21 (5) (2019) 1235–1247.

[51] VQEG, Final report from the video quality experts group on the validation of objective models of video quality assessment, Tech. rep., VQEQ (2000).

[52] P. Hanhart, T. Ebrahimi, Calculation of average coding efficiency based on subjective quality scores, Journal of Visual Communication and Image Representation 25 (2014) 555–564.