# Image Quality Assessment in the Modern Age

Kede Ma
City University of Hong Kong
Hong Kong SAR, China
kede.ma@cityu.edu.hk

Yuming Fang
Jiangxi University of Finance and Economics
Jiangxi, China
fa0001ng@e.ntu.edu.sg

## ABSTRACT

This tutorial provides the audience with the basic theories, methodologies, and current progresses of image quality assessment (IQA). From an actionable perspective, we will first revisit several subjective quality assessment methodologies, with emphasis on how to properly select visual stimuli. We will then present in detail the design principles of objective quality assessment models, supplemented by an in-depth analysis of their advantages and disadvantages. Both hand-engineered and (deep) learning-based methods will be covered. Moreover, the limitations with the conventional model comparison methodology for objective quality models will be pointed out, and novel comparison methodologies such as those based on the theory of "analysis by synthesis" will be introduced. We will last discuss the real-world multimedia applications of IQA, and give a list of open challenging problems, in the hope of encouraging more and more talented researchers and engineers devoting to this exciting and rewarding research field.

## 1 INTRODUCTION

Image quality assessment (IQA), a long-standing task in the field of image and multimedia processing, has evolved rapidly in the past two decades [23], and has also gained increasing attention from both academic and industry for its broad applications. In this extended abstract, we plan to divide and introduce IQA in the following four parts:

- Subjective IQA, the most straightforward and reliable way of assessing perceptual quality by humans;
- Objective IQA, constructing computational models to automate the quality assessment process;
- IQA model comparison, quantifying the (relative) quality prediction performance of the competing models;
- IQA model applications, considering the particularities of different forms of multimedia data.

## 2 SUBJECTIVE IQA

The goal of subjective IQA is to collect *reliable* mean opinion scores (MOSs) from human subjects on the perceived quality of test images. Several subjective methodologies have been standardized in the ITU-R and ITU-T recommendations [2], which can be broadly categorized into single-stimulus, double-stimulus, and multiple-stimulus methods. Take the single-stimulus absolute category rating (ACR) as an example. Each test image is rated individually using the labels "bad", "poor", "fair", "good", and "excellent", which are translated to the values 1, 2, 3, 4, and 5 when calculating the MOS. Along with the introduction of subjective experimental procedures, many important (but subtle) designs are also discussed.

- Which subset of test images are "ideal" to choose from the web-scale unlabeled database for human annotation?
- How much instruction should be provided to subjects for more consistent and less biased MOS collection?
- What are the general guidelines to set up the experimental environment, especially the viewing conditions?

The immediate results of subjective experiments are human-labeled image quality databases, which monitor the progress of objective IQA. For example, the LIVE dataset [18] marks the switch from distortion-specific to general-purpose IQA. The CSIQ dataset [10] enables cross-dataset comparison. The TID2013 dataset [17] and its successor KADID-10K [11] expose the difficulty of IQA methods in generalizing to different distortion types. The Waterloo Exploration Database [14] tests model robustness to diverse content variations of natural scenes. The LIVE Challenge Database [7] probes the synthetic-to-real generalization, which is further evaluated by the KonIQ-10K [8] and SPAQ [6] datasets. At the end of this part, we will give a brief overview of these datasets, and share our thoughts on creating better IQA databases in terms of mining hard and diverse images and collecting reliable MOSs.

## 3 OBJECTIVE IQA

Objective IQA aims to develop computational algorithms that are capable of providing consistent quality predictions with human data. These models can be mainly classified into two categories: full-reference (FR) and no-reference (NR, or blind) models. FR-IQA methods assume full access to a pristine undistorted image (also referred to as the reference image) for quality assessment of a "distorted" image. NR-IQA models, on the other hand, do not require any reference information. We will first discuss full-reference models, and start with THE default quality metric - mean squared error (MSE) that has dominated the field of signal processing for more than 50 years [24]. We will revisit the limitations of MSE by hand-crafting its counterexamples intuitively. This motivates the

development of the structural similarity (SSIM) index [22], a award-winning and widely adopted perceptual quality model. Since its inception in 2004, the design philosophy underlying SSIM continues to impact the IQA field up to today. Among a myriad of existing IQA models, we will sample a few that we believe advance the field from at least one of the following aspects:

- More accurate IQA in terms of explaining human data in existing databases;
- Color IQA that gives a better account for color perception of the human visual system;
- Misalignment-aware IQA that does not require the reference and distorted images to be precisely aligned;
- Texture-aware IQA that provides an efficient characterization of texture similarity;
- IQA based on other design principles, e.g., information theoretic and data-driven approaches.

We will conclude the discussion of FR-IQA models by pointing out an embarrassing and common design flaw: many IQA models fail to satisfy the identity of indiscernibles[1], which has a strong implication that they are not suitable for perceptual optimization.

We then switch our attention to NR-IQA, which is more practical and challenging due to the lack of reference information. We will first describe a widely accepted design principle based on natural scene statistics (NSS) [25]. The underlying assumption is that a measure of the destruction of statistical regularities of natural images provides a reasonable approximation to perceived visual quality. Both hand-crafted and learned NSS in spatial and frequency domain will be described. In particular, we would like to put more emphasis on one NR model, namely, the naturalness image quality evaluator (NIQE) [16], which has began to show its potentials in benchmarking image processing algorithms in real settings.

Limited by the expressiveness of hand-crafted features, NSS-based approaches have been surpassed by data-driven NR-IQA models based on convolutional neural networks in recent years. Patch-wise training, transfer learning, and quality-aware pre-training are means of compensating for the lack of human data. Apart from summarizing the specialized architectural designs, we plan to draw the audience's attention to the latest learning paradigms for NR-IQA, including

- Unified learning for NR-IQA from multiple IQA databases simultaneously without additional subjective testing for perceptual scale realignment [29];
- Active learning for NR-IQA by failure identification and model rectification [27];
- Continual learning for NR-IQA, where the model evolves with new data while being resistant to catastrophic forgetting of old data [28].

## 4 IQA MODEL COMPARISON

Conventional IQA model comparison generally follows a three-step approach. First, pre-select a number of images to form the test set. Second, collect the MOS for each image in the test set to represent its true perceptual quality. Third, rank the competing models

according to their goodness of fit (e.g., Spearman rank-order correlation coefficient) on the test set. The one with the best result is declared the winner. We will discuss the limitations of this conventional method in terms of the representativeness of test samples and the risk of overfitting. We will then introduce a series of alternative IQA model comparison methods, including

- Maximum differentiation (MAD) competition [26], automatically synthesizing images that are likely to falsify the IQA model in question;
- Group MAD (gMAD) competition [13], a discrete instantiation of MAD that is more efficient and controllable;
- Eigen-distortion analysis [1], a method for comparing image representations in terms of their ability to explain perceptual sensitivity in humans;
- Comparison of IQA models for perceptual optimization of image processing systems [3].

All the above-mentioned methods are based on the idea of "analysis by synthesis", which is rooted in the pattern theory by Ulf Grenander.

## 5 IQA MODEL APPLICATIONS

It is highly nontrivial to apply IQA techniques in the field of multimedia due to substantially different data formats and particularities [21]. Subject to the time constraint, we plan to present a few demonstrating examples, including

- High-dynamic-range imaging [9], where the input and output images have different bit depths;
- Image fusion [12], where input and output have different numbers of images;
- Color-to-gray conversion [15] and colorization, where input and output images have different color channels;
- Image retargeting [5], where input and output images have different spatial resolutions;
- Stereoscopic images [20], where binocular vision should be modeled;
- Omnidirectional images [19], where viewing behaviors may be indispensable for quality assessment;
- Screen content images [4], where non-natural image statistics should be extracted;
- Natural videos (in the streaming setting), where the time dimension is added, leading to complex spatiotemporal distortions.

We will definitely point to the audience useful resources for IQA applications that have respectfully not covered. We will also cover some general and intuitive applications of IQA such as automatic hyperparameter adjustment and optimization of image processing algorithms.

As a final remark, through this tutorial, we sincerely hope more and more talented researchers and engineers are willing to join us, contributing to this exciting and rewarding field.

## REFERENCES

[1] A. Berardino, J. Ballé, V. Laparra, and E. P. Simoncelli. 2017. Eigen-distortions of hierarchical representations. *arXiv preprint arXiv:1710.02266* (2017).
[2] RECOMMENDATION ITU-R BT. 2002. Methodology for the subjective assessment of the quality of television pictures. *International Telecommunication Union* (2002).

---

[1]Give an FR-IQA model $D(\cdot)$, where a lower score indicating better predicted quality with a minimum of zero, and two images $x$, $y$, the identity of indiscernibles refers to $D(x, y) = 0 \Leftrightarrow x = y$.

[3] K. Ding, K. Ma, S. Wang, and E. P. Simoncelli. 2021. Comparison of full-reference image quality assessment models for optimization of image processing systems. *International Journal of Computer Vision* 129, 4 (2021), 1258–1281.

[4] Y. Fang, J. Yan, J. Liu, S. Wang, Q. Li, and Z. Guo. 2017. Objective quality assessment of screen content images by uncertainty weighting. *IEEE Transactions on Image Processing* 26, 4 (2017), 2016–2027.

[5] Y. Fang, K. Zeng, Z. Wang, W. Lin, Z. Fang, and C.-W. Lin. 2014. Objective quality assessment for image retargeting based on structural similarity. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems* 4, 1 (2014), 95–105.

[6] Y. Fang, H. Zhu, K. Ma, Z. Wang, and S. Li. 2020. Perceptual evaluation for multi-exposure image fusion of dynamic scenes. *IEEE Transactions on Image Processing* 29 (2020), 1127–1138.

[7] D. Ghadiyaram and A. C. Bovik. 2015. Massive online crowdsourced study of subjective and objective picture quality. *IEEE Transactions on Image Processing* 25, 1 (2015), 372–387.

[8] V. Hosu, H. Lin, T. Sziranyi, and D. Saupe. 2020. KonIQ-10k: An ecologically valid database for deep learning of blind image quality assessment. *IEEE Transactions on Image Processing* 29 (2020), 4041–4056.

[9] V. Laparra, J. Ballé, A. Berardino, and E. P. Simoncelli. 2016. Perceptual image quality assessment using a normalized Laplacian pyramid. *Electronic Imaging* 2016, 16 (2016), 1–6.

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

[11] H. Lin, V. Hosu, and D. Saupe. 2019. KADID-10k: A large-scale artificially distorted IQA database. In *2019 Eleventh International Conference on Quality of Multimedia Experience*. 1–3.

[12] Z. Liu, E. Blasch, Z. Xue, J. Zhao, R. Laganiere, and W. Wu. 2011. Objective assessment of multiresolution image fusion algorithms for context enhancement in night vision: A comparative study. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 34, 1 (2011), 94–109.

[13] K. Ma, Z. Duanmu, Z. Wang, Q. Wu, W. Liu, H. Yong, H. Li, and L. Zhang. 2020. Group maximum differentiation competition: Model comparison with few samples. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 42, 4 (2020), 851–864.

[14] K. Ma, Z. Duanmu, Q. Wu, Z. Wang, H. Yong, H. Li, and L. Zhang. 2017. Waterloo exploration database: New challenges for image quality assessment models. *IEEE Transactions on Image Processing* 26, 2 (2017), 1004–1016.

[15] K. Ma, T. Zhao, K. Zeng, and Z. Wang. 2015. Objective quality assessment for color-to-gray image conversion. *IEEE Transactions on Image Processing* 24, 12 (2015), 4673–4685.

[16] A. Mittal, R. Soundararajan, and A. C. Bovik. 2013. Making a "completely blind" image quality analyzer. *IEEE Signal Processing Letters* 20, 3 (2013), 209–212.

[17] N. Ponomarenko, O. Ieremeiev, V. Lukin, K. Egiazarian, L. Jin, J. Astola, B. Vozel, K. Chehdi, M. Carli, F. Battisti, et al. 2013. Color image database TID2013: Peculiarities and preliminary results. In *European Workshop on Visual Information Processing*. 106–111.

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

[19] X. Sui, K. Ma, Y. Yao, and Y. Fang. 2021. Perceptual quality assessment of omnidirectional images as moving camera videos. *IEEE Transactions on Visualization and Computer Graphics* (2021), 1–11.

[20] J. Wang, A. Rehman, K. Zeng, S. Wang, and Z. Wang. 2015. Quality prediction of asymmetrically distorted stereoscopic 3D images. *IEEE Transactions on Image Processing* 24, 11 (2015), 3400–3414.

[21] Z. Wang. 2016. Objective image quality assessment: Facing the real-world challenges. *Electronic Imaging* 2016, 13 (2016), 1–6.

[22] Z. Wang, A.C. Bovik, H.R. Sheikh, and E. P. Simoncelli. 2004. Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing* 13, 4 (2004), 600–612.

[23] Z. Wang and A. C. Bovik. 2006. *Modern Image Quality Assessment.* Morgan & Claypool.

[24] Z. Wang and A. C. Bovik. 2009. Mean squared error: Love it or leave it? A new look at signal fidelity measures. *IEEE Signal Processing Magazine* 26, 1 (2009), 98–117.

[25] Z. Wang and A. C. Bovik. 2011. Reduced-and no-reference image quality assessment. *IEEE Signal Processing Magazine* 28, 6 (2011), 29–40.

[26] Z. Wang and E. P. Simoncelli. 2008. Maximum differentiation (MAD) competition: A methodology for comparing computational models of perceptual quantities. *Journal of Vision* 8, 12 (2008), 8.1–8.13.

[27] Z. Wang, H. Wang, T. Chen, Z. Wang, and K. Ma. 2021. Troubleshooting blind image quality models in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition*. 16256–16265.

[28] W. Zhang, D. Li, C. Ma, G. Zhai, X. Yang, and K. Ma. 2021. Continual learning for blind image quality assessment. *arXiv preprint arXiv:2102.09717* (2021).

[29] W. Zhang, K. Ma, G. Zhai, and X. Yang. 2021. Uncertainty-aware blind image quality assessment in the laboratory and wild. *IEEE Transactions on Image Processing* 30 (2021), 3474–3486.