

BVI-UGC: A Video Quality Database for User-Generated Content Transcoding

Zihao Qi, *Student Member, IEEE*, Chen Feng, *Student Member, IEEE*, Fan Zhang, *Member, IEEE*, Xiaozhong Xu, *Senior Member, IEEE*, Shan Liu, *Fellow, IEEE*, and David R. Bull, *Fellow, IEEE*

Abstract—In recent years, user-generated content (UGC) has become one of the major video types consumed via streaming networks. Numerous research contributions have focused on assessing its visual quality through subjective tests and objective modeling. In most cases, objective assessments are based on a no-reference scenario, where the corresponding reference content is assumed not to be available. However, full-reference video quality assessment is also important for UGC in the delivery pipeline, particularly associated with the video transcoding process. In this context, we present a new UGC video quality database, BVI-UGC, for user-generated content transcoding, which contains 60 (non-pristine) reference videos and 1,080 test sequences. In this work, we simulated the creation of non-pristine reference sequences (with a wide range of compression distortions), typical of content uploaded to UGC platforms for transcoding. A comprehensive crowdsourced subjective study was then conducted involving more than 3,500 human participants. Based on this collected subjective data, we benchmarked the performance of 10 full-reference and 11 no-reference quality metrics. Our results demonstrate the poor performance (SROCC values are lower than 0.6) of these metrics in predicting the perceptual quality of UGC in two different scenarios (with or without a reference). To facilitate future research in this area, we have made BVI-UGC publicly available at <https://zihaoq1.github.io/BVI-UGC/>

Index Terms—Video quality assessment, UGC, video transcoding, BVI-UGC, subjective study, crowdsourcing

I. INTRODUCTION

With advances in mobile devices and communication network technologies, coupled with the explosion of social media and streaming platforms, user-generated content (UGC) now represents more than 35% of downstream volume over fixed networks [1], streamed by service providers such as YouTube, Facebook, TikTok, and Tencent. UGC also significantly influences upstream traffic, driven by the popularity of short-form videos [1]. These highlight the importance of video compression in UGC streaming, which plays a critical role in managing the trade-offs between video quality and required bandwidth.

Compared to professionally-generated content (PGC), UGC has unique characteristics due to the specific production and

Zihao Qi, Chen Feng, Fan Zhang and David Bull are with the Visual Information Lab, University of Bristol, Bristol BS1 5DD, U.K. (e-mail: {zihao.qi, chen.feng, fan.zhang, dave.bull@bristol.ac.uk}). Xiaozhong Xu and Shan Liu are with Tencent America, Palo Alto, USA. (e-mail: {xiaozhongxu, shanliu}@tencent.com)

The authors acknowledge the funding from Tencent (US), University of Bristol, and the UKRI MyWorld Strength in Places Programme (SIPF00006/1).

This work involved collecting data from human participants. The relevant experiments have been approved by the Faculty of Engineering Research Ethics Committee of the University of Bristol (Ref 12352).

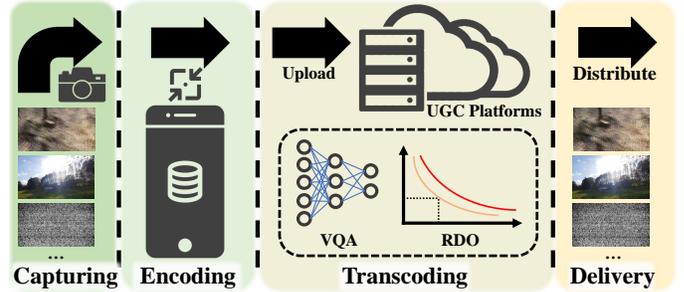


Fig. 1: Illustration of the UGC video delivery pipeline. Source videos captured by users may contain various distortions due to poor quality equipment, unskilled cinematography and lossy compression. Captured videos are uploaded to UGC platforms where they are transcoded and streamed to video consumers. During the transcoding process, a perceptually accurate VQA metric is a key component in rate-distortion optimization.

delivery pipeline employed (Fig. 1). UGC videos are typically captured using commercial and mobile devices by amateur users. In most cases, the original uncompressed sources of these videos are not stored during acquisition, as they are directly compressed using a fast video codec embedded in the device (e.g., x264 [2]). The compressed videos are then uploaded to a selected UGC platform for streaming, where a further layer of compression (i.e., transcoding) is performed. Due to the nature of UGC, its transcoding is often based on references, which themselves contain significant source and compression artifacts. This differs from the PGC production pipeline, where high-quality original content is always available during encoding.

In current transcoding pipelines used by many UGC platforms, the encoding operation is similar to that used for PGC, where the objective is to minimize the distortion (or the visual quality degradation) between the reference and reconstructed videos. As mentioned above, in UGC transcoding, the input reference content may contain various source and compression artifacts, which reduces the effectiveness of the rate-distortion (or rate-quality) optimization due to the inaccurate quality (distortion) prediction in these cases [3].

In recent years, there have been many research contributions addressing the issues associated with UGC-based video quality assessment. However, most solutions [3–11] have been developed for no-reference scenarios, which are typically employed to assess visual quality at the user end. Although there are existing studies that address the subjective quality assessment of transcoded UGC, most of these only consider a limited

range of reference quality levels [12–14], which do not reflect real-world scenarios where the quality of reference content varies considerably.

In this context, to facilitate advances in video quality assessment for UGC, we have developed a large-scale UGC database, BVI-UGC, focused on transcoding applications. We collected 60 high-quality source sequences, covering 15 major UGC categories. These videos were then compressed using the x264 codec [2] with three quantization levels that simulate the compression operation during content capture. The ingested content (i.e., non-pristine reference used in the transcoding stage) was further transcoded using three video codecs (x264, x265 [15] and libaom [16]) with two resolution re-sampling levels and three quantization parameters, resulting in a total number of 1,080 distorted test sequences. Moreover, an extensive crowdsourcing-based subjective experiment was performed to collect quality scores from more than 3,500 participants for both test sequences and non-pristine references. The ground-truth data has been further employed to benchmark 21 existing full-reference (FR) and no-reference (NR) objective quality models. The experimental results highlight the poor and inconsistent performance of existing metrics and confirm the urgent requirement for more accurate quality assessment methods. The primary contributions of this work are summarized below.

- 1) BVI-UGC is the **first public UGC video quality databases containing reference clips compressed with various quantization levels**, which simulates the UGC delivery pipeline. All the source, non-pristine references and test sequences are made available for public evaluation.
- 2) All test and reference sequences are labeled with ground truth subjective quality scores through a **reliable, large-scale crowdsourcing-based psychophysical experiment** using Amazon Mechanical Turk [17] platform. We also open-sourced the web application used for this subjective test.
- 3) We demonstrate how to exploit this database by **benchmarking FR and NR VQA methods in the context of UGC transcoding**. We have also employed this database to evaluate the performance of NR quality metrics in terms of directly measuring the perceptual quality of test sequences without references. Based on a comprehensive experiment involving 21 full-/no-reference metrics, we highlight the challenging nature (no existing quality metrics achieve a SROCC value higher than 0.6 on this dataset) of UGC-based video quality assessment in the context of transcoding applications.

The remainder of this paper is organized as follows. Section II provides a brief summary of related work in the research areas of UGC video quality assessment and existing UGC databases. Section III describes the development of the database, including source video collection and capture, and reference/test sequence generation. Section IV presents the detailed methodology and configuration employed in the subjective test, and performs an analysis of the collected subjective quality scores. Section V summarizes the results of the benchmark experiment for 21 video quality metrics.

Finally, Section VI provides a conclusion of the paper and outlines future work.

II. RELATED WORKS

In this section, we first review previous work addressing full-reference and no-reference video quality assessment (VQA), in particular those developed specifically for UGC. We then summarize existing UGC video quality databases, and highlight the urgent need to develop a more diverse video quality database for UGC transcoding.

A. Video Quality Assessment for UGC

Although subjective tests provide a gold standard for estimating the perceptual quality of video content, they are not widely employed in practical applications due to their time-consuming, non-real-time and expensive nature [18]. Instead, objective video quality assessment methods are frequently used in algorithm benchmarking and optimization. These methods can be classified according to the availability of reference content, into two major categories¹: full-reference (FR) models that provide a quality prediction based on the comparison between a processed sequence and the reference counterpart, and no-reference (NR) models, which directly assess the quality of a sequence without considering any information from its reference content.

1) *Full-Reference VQA*: FR VQA models typically measure the difference between a test video sequence and its corresponding reference. Simple models such as PSNR are widely employed across many image and video processing applications, serving as a benchmark metric for algorithm comparison and in loss functions for model optimization (e.g., rate-distortion optimization in compression). To further improve their correlation performance with perceptual quality, researchers have developed perceptually-inspired quality metrics that exploit various characteristics of the human vision system (HVS) such as texture masking [19], just noticeable difference [20] and contrast sensitivity functions [21]. Notable examples include SSIM and its variants [22–27], ADM [28], VIF [29], MOVIE [30], MAD [31] and PVM [32]. Moreover, some of these perceptual models have been combined together with various video features through linear regression (e.g. SVM [33]) to achieve even higher correlation performance, with one of the most successful examples being VMAF [34].

More recently, researchers have focused on deep learning-based solutions for quality assessment. Important contributions in this category include DeepVQA [35], LPIPS [36], C3DVQA [37], and CUGCVQA [38]. Although these methods show promise when compared to conventional and regression-based methods, in most cases they demand an intra-database cross-validation due to their limited model generalization ability. To address this issue, more recently, unsupervised or weakly-supervised learning strategies have been developed for deep VQA models [39–41], which do not require training data with ground-truth subjective scores.

¹There is another class of objective VQA method, denoted reduced-reference (RR), used when only partial information from the reference is available.

TABLE I: Features of notable user-generated content video databases with transcoded sequences. Here ‘‘L’’ stands for ‘Landscape’ layout, while ‘‘P’’ stands for ‘Portrait’. [†]BVI-UGC contains source videos collected from YouTube-UGC database and captured by lab participants using various devices.

	YT-UGC VP9	LIVE-WILD	UGC-VIDEO	ICME 2021	TaoLive	BVI-UGC
source seq.	n/a	n/a	n/a	n/a	n/a	60
reference seq.	169	55	50	1000	418	60
transcoded seq.	567	220	500	7000	3,344	1080
ref. quality	low	medium	medium	high	high	3 levels
codecs	VP9	x264	x264 and x265	x264	x265	x264,x265 and libaom
content source	YouTube	Mobile captured	Mobile captured	Mobile captured	TaoLive	Mixed [†]
resolution	720p,1080p	360p,540p,720p,1080p	720p	720p	720p,1080p	540p,1080p
layouts	L	L	P	L & P	mostly P	L & P
frame rate	30fps	24-30fps	24-30fps	30fps	20-30fps	24-60fps
bit depth	8	8	8	8	8	8
duration	20s	10s	10s	5s	8s	5s
rating scale	Continuous 1-5	0-100	Discrete 1-5	Discrete 1-5	Discrete 1-5	0-100
subject number	n/a	40	28	n/a	44	3,500+
ratings Avg.	n/a	20	28	>50	44	160

When employed in the context of UGC transcoding, where reference content is often non-pristine, containing various visible artifacts, the aforementioned FR VQA methods do not offer satisfactory correlation performance with ground-truth subjective data [3, 42]. This can lead to inconsistent quality prediction when used for algorithm evaluation and poor rate quality optimization performance if employed in the transcoding loop.

2) *No-reference VQA*: NR quality metrics are employed when the reference content is not available during quality assessment. In video delivery, this is typically applied at the decoder (user end) to measure the quality of user experience. Numerous NR quality assessment methods have been developed for estimating compression distortion [43, 44], transmission errors [45, 46] and specific artifacts [47, 48].

Similar to FR VQA, recent NR quality metrics have exploited deep neural networks to enhance quality assessment. These have demonstrated improved performance over conventional NR VQA methods based on classical signal processing theories. These learning-based methods can be classified as supervised or weakly-/un-supervised approaches. The former directly performs model learning based on ground-truth subjective scores collected from human participants in large scale experiments, with notable examples including V-MEON [49], ChipQA [50], Compressed VQA [38] SimpleVQA [51], FastVQA [10] and FasterVQA [11]. Un-supervised or weakly-supervised methods typically convert the main task (directly predicting quality scores) to an auxiliary mission based on techniques such as contrastive learning or ranking learning. This supports the generation of more diverse training content without performing expensive subjective tests. Important work in this class includes NR-RankDVQA [39], CONTRIQUE [40] and CONVIQT [52].

Many NR quality metrics focus on assessing the quality of the UGC content, such as [4, 38, 40, 51, 52]. However, their performance has not been fully investigated in the context of UGC transcoding due to the limited availability of benchmark databases of this type.

B. UGC Video Quality Databases

To evaluate the performance of objective quality metrics, subjective video quality databases that contain various distorted (and reference, if for full-reference scenarios) sequences with different visual artifacts, are used. Human participants are employed to view these sequences through psychophysical experiments based on certain test methodologies and procedures [53]. Different correlation coefficients between the quality indices generated by the objective quality models for these test sequences and their corresponding ground-truth subjective scores, collected in the subjective experiment, can then be used to measure and compare the performance of these objective quality models.

Early work on subjective video databases typically focused on PGC videos, with notable examples such as VQEG FR Phase I, VQEG HD [54], LIVE-VQA [55], IVP [56] and BVI-HD [57], which primarily investigate the impact of video compression and transmission. There are also contributions that study the influence of specific video formats or artifacts, including BVI-HFR (frame rate) [58], LIVE-YT-HFR (frame rate) [59, 60], BVI-SR (spatial resolution) [61], BVI-BD (bit depth) [62], BAND-2k (banding artifacts) [63, 64], etc.

In the context of UGC, most subjective quality assessment approaches focus on the no-reference scenario. Existing publicly available UGC databases such as YouTube UGC [3], KoNViD-1k [5, 6], KoNViD-150k [7, 8], and LIVE-VQC [9] only provide distorted (transcoded) sequences in the absence of their corresponding (non-pristine) reference or (pristine) original content. Only a few databases attempt to extend research to the full-reference (transcoded) case. For example, YouTube-UGC [3] contains a subset that provides VP9 transcoded content together with their low-quality references. The Quality Assessment Grand Challenge of ICME 2021 [14] released a large UGC database containing 7000 transcoded sequences and their 1000 nearly-pristine references. Similarly, TaoLive [65] is another large-scale database that provides high-quality references together with transcoded content. LIVE-WILD [13] and UGC-VIDEO [12] are also notable works in this area, consisting of test sequences compressed from medium quality level references. The main features of

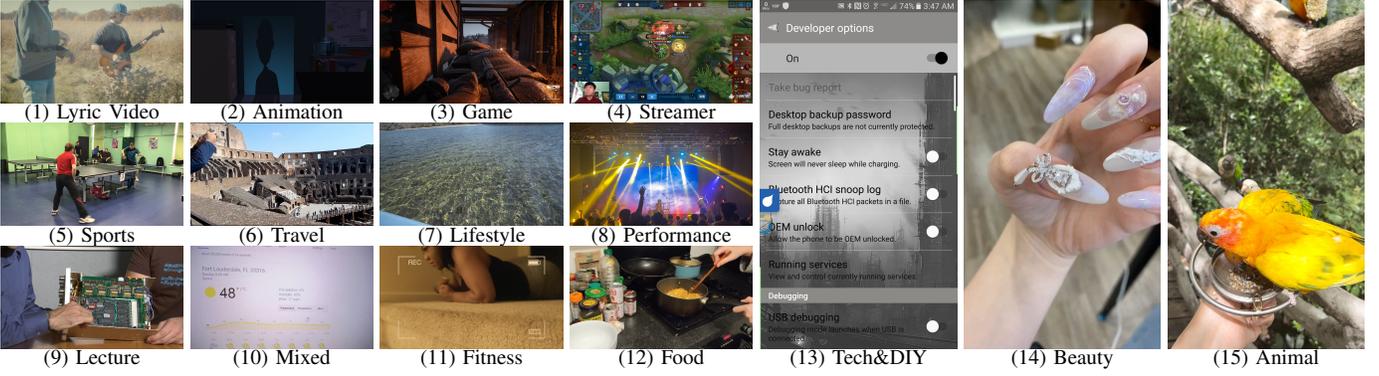


Fig. 2: Example frames in selected videos from 15 UGC categories defined in this database. Every category contains both landscape and portrait videos (the ratio is approximately 2:1).

these four databases are summarized in TABLE I. It is noted that none of these contains references at multiple quality levels, and their transcoding process is often based on a single video codec (except for the UGC-VIDEO database), which does not reflect the diversity of the transcoded UGC videos.

III. THE BVI-UGC DATABASE

To address the issues mentioned in Section II regarding the lack of comprehensive video quality databases for UGC transcoding, we have developed a new database, BVI-UGC which, for the first time, specifically considers the UGC transcoding process. We simulated the UGC video delivery pipeline shown in Fig. 1 to generate non-pristine references at different quality levels, and applied another layer of encoding (transcoding) using three commonly used codecs to obtain distorted video sequences.

This section describes the methodology used to select/capture source sequences in the BVI-UGC database, and the workflow to generate non-pristine reference and transcoded content.

A. Source Sequences

In order to develop a database with diverse and representative source content, we first defined 15 typical UGC categories following the scope of existing UGC databases and the genres on popular UGC platforms (e.g. YouTube, Tiktok, Flickr etc.) [3], including (1) Lyric Video; (2) Animation; (3) Game; (4) Streamer; (5) Sports; (6) Travel; (7) Lifestyle; (8) Performance; (9) Lecture; (10) Mixed; (11) Fitness; (12) Food; (13) Tech & DIY; (14) Beauty; (15) Animal. In each category, ten videos were collected from the YouTube-UGC database [3] or captured using mobile phones and drones. The scenes captured in these sequences were designed to diversify low and high level features, such as spatial textures, motions, lighting conditions, backgrounds and foregrounds. This results in a total of 150 ten-second candidate videos. It should be noted that we ensured that all the source sequences are visually lossless so that they can be used as ‘pristine’ original content in the UGC pipeline.

The 150 candidate videos were then truncated into 300 short (5 second) clips, following the optimal duration study conducted in [66, 67]. This allows us to generate more test

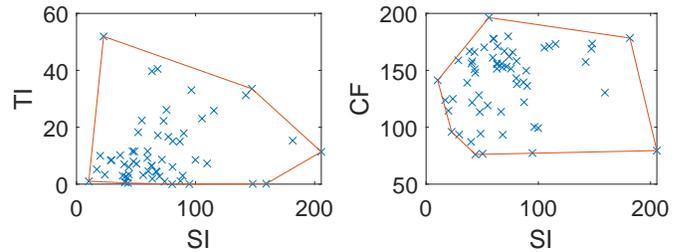


Fig. 3: Feature distribution of the source content in the BVI-UGC database. (Left) SI versus TI; (right) SI versus CF.

sequences given the limited time and financial resources. To support further content selection, we follow the procedures in [57, 68] to determine 60 final source sequences (four for each UGC category). It is also noted that, considering that many UGC videos are associated with a portrait layout, we selected both landscape and portrait videos with a ratio (between landscape and portrait content) of 2:1. Example frames of the selected videos from each of the 15 categories are shown in Fig. 2. To showcase the content distribution of the collected content, we calculated these primary video features for each source sequence, including Spatial Information (SI), Temporal Information (TI) and Colorfulness (CF) for each source sequence, following the feature definitions in [69], with the average values of these features (at the sequence level) shown in Fig. 3. It can be observed that the collected content covers a relatively wide range for each video feature compared to other databases in the literature [4, 57, 70].

B. Non-pristine References

In a real-world UGC production pipeline, pristine original content is not available on a user’s device due to its large storage requirement. Instead, it is typically compressed using a fast video codec to generate compressed sequences that may contain visual artifacts. Their quality may vary, as different coding configurations can be employed here. To simulate this process, we compressed 60 source sequences using a fast implementation of H.264/AVC [71], x264 [2], based on its *slow* preset, which is one of the most commonly used video codecs in practical applications. To further diversify the quality of reference sequences, three quantization parameters were used to generate content with high (QP=30), medium (QP=37)

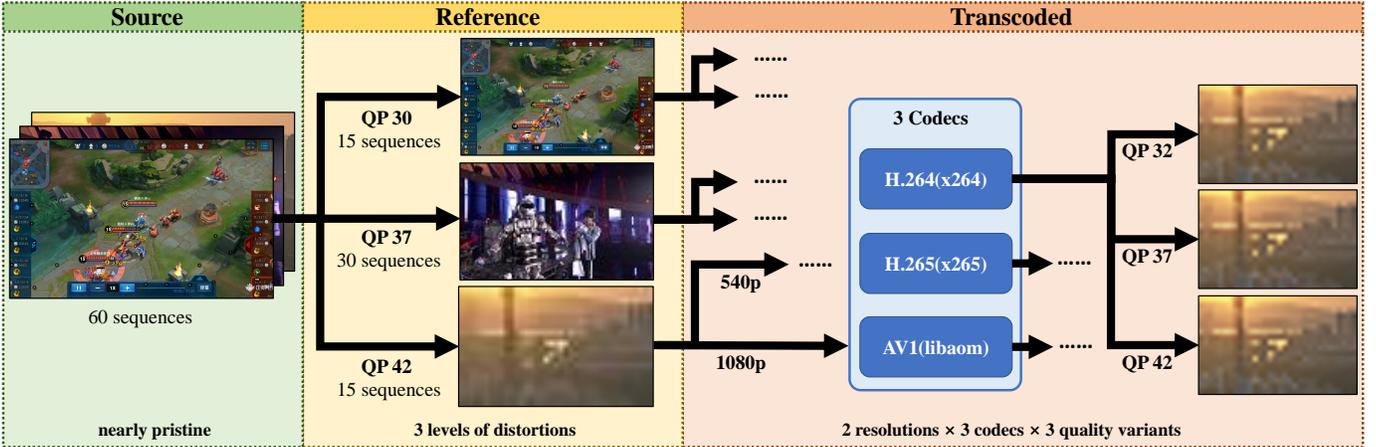


Fig. 4: Illustration of the content generation process for the BVI-UGC database, which contains 60 non-pristine reference and 1080 transcoded sequences.



Fig. 5: Sample blocks from the high quality source, non-pristine reference and transcoded videos (for the same content).

and low (QP=42) visual quality levels, as shown in Fig. 4. For each source sequence, only one QP value was employed for compression, and we ensured that there is at least one video in every content category encoded with each QP level. This results in 15 non-pristine reference sequences in the QP30 group, 30 in the QP37 group, and 15 in the QP42 group. We generated more non-pristine reference sequences in the medium-quality QP group based on the assumption that this is more common in the real-world UGC production scenarios.

C. Distorted Sequences

Sixty non-pristine references are then compressed again to simulate the transcoding operation. In order to generate more diverse content, we used three video codecs, x264 (ffmpeg v4.4.1 built-in [72], *slow* preset), x265 (ffmpeg v4.4.1 built-in [72], *slow* preset), and libaom [73] (v1.0.0, Random Access [74]), which are commonly employed by various UGC streaming platforms. For each codec, three different QP values are used to produce content at various quality levels. They are {32, 37 and 42} for x264 and x265, and {43, 55 and 63} for libaom. To further emulate the real-world streaming

scenarios, resolution adaptation has also been applied (with a factor of two) in compression. This results in a total number of 1,080 transcoded sequences (60 references \times 2 resolutions \times 3 codecs \times 3 quality variants, 18 per reference). The content generation workflow is illustrated in Fig. 4.

Visual examples have been provided to demonstrate quality differences between high-quality sources, distorted references, and transcoded content, as shown in Fig. 5. It can be observed that, when the distorted reference is of low quality (i.e. compressed by x264 using QP42), it is likely that the transcoded sequence will be associated with slightly better perceptual quality due to the artifact filtering and smoothing effect through compression; for example, in Fig. 5, the transcoded content by x265 (QP37) looks better than the unpristine reference sequence (x264, QP42). This could challenge FR VQA methods, most of which assume perfect quality with reference content.

IV. SUBJECTIVE EXPERIMENTS

This section describes the configuration of the crowdsourcing-based subjective experiment, which collected the quality opinion scores of the video sequences in the BVI-UGC database. We then analyze the subjective results and demonstrate its reliability.

A. Experiment Design

Due to the large number of video sequences in this database and the nature of UGC consumption (typically within inconsistent viewing conditions and on different devices), we designed a crowdsourcing subjective test based on the Amazon Mechanical Turk [17] platform. In this experiment, we employed the Absolute Category Rating with Hidden Reference (ACR-HR) methodology [75], which has been commonly used for many subjective UGC studies [3, 12–14], to collect subjective scores for 60 non-pristine reference and 1,080 transcoded sequences. High-quality source content was not shown in this experiment, but is provided alongside other sequences in the database.

In each test session, prior to the formal test, subjects were first asked to calibrate their screen resolution based on the

Algorithm 1: The calculation of MOS based on [53].

Input: Raw subjective scores $\{o_{ij}\}$, subject index $i \in \{1, \dots, I_j\}$ and sequence index $j \in \{1, \dots, J_i\}$, in which I_j are J_i the number of subjects that have scored sequence j and the number of sequences scored by subject i , respectively.

Output: Weighted mean opinion score MOS_j , standard error of score SE_j , subject inconsistency σ_i and subject bias b_i

```

1 Initialization
2  $MOS_{j,0} \leftarrow \frac{1}{I_j} \sum_{i=1}^{I_j} o_{ij}$ 
3  $b_i \leftarrow \frac{1}{J_i} \sum_{j=1}^{J_i} (o_{ij} - MOS_{j,0})$ 
4  $t \leftarrow 0$ 
5 while  $t < 1000$  do
6    $r_{ij} = o_{ij} - MOS_{j,t} - b_i$ 
7    $r_i = \frac{1}{J_i} \sum_{j=1}^{J_i} r_{ij}$ 
8    $\sigma_i = \sqrt{\frac{1}{J_i} \sum_{j=1}^{J_i} (r_{ij} - r_i)^2}$ 
9    $MOS_{j,t+1} = \sum_{i=1}^{I_j} \sigma_i^{-2} (o_{ij} - b_i) / \sum_{i=1}^{I_j} \sigma_i^{-2}$ 
10   $b_i = \frac{1}{J_i} \sum_{j=1}^{J_i} (o_{ij} - MOS_{j,t+1})$ 
11  if  $\sum_{j=1}^{J_i} (MOS_{j,t+1} - MOS_{j,t})^2 < 10^{-16}$  then
12    break
13  end
14   $t = t + 1$ 
15 end
16  $r_j = \frac{1}{I_j} \sum_{i=1}^{I_j} r_{ij}$ 
17  $SE_j = \frac{1}{I_j} \sqrt{\sum_{i=1}^{I_j} (r_{ij} - r_j)^2}$ 

```

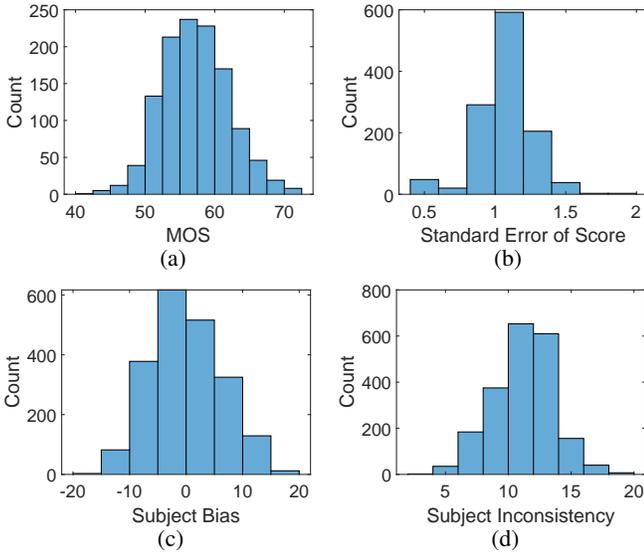


Fig. 6: Histogram of (a) the MOS; (b) standard error of score over subjects; (c) subject bias; (d) subject inconsistency.

method developed in [76] and then perform a vision acuity test using Ishihara and Snellen charts. If participants fail in this test, the session is stopped without showing any video sequences. For those who have passed the test, they are shown video sequences randomly picked from 60 non-pristine reference and 1,080 transcoded videos, each of which is played only once. After viewing each video, participants are asked to provide a subjective score of video quality using a continuous slider with five evenly spaced intervals labeled *Bad*, *Poor*, *Fair*, *Good* and *Excellent*, each of which covers a quality range

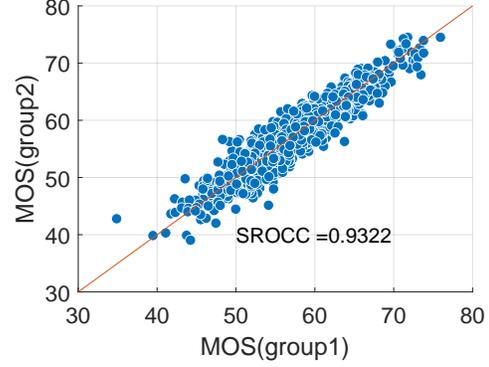


Fig. 7: Scatter plot of MOS values given by two randomly separated halves of subjective data, from one random split in the 1,000 repetitions. Higher correlated ratio indicates higher inter-subject consistency.

of 20 respectively. To allow subjects to familiarize with the experiment, 3 training trials are presented before the formal test. During the test, participants are free to leave anytime, but the results collected from participants who viewed less than 20 videos were discarded to make sure sufficient videos are watched by each subject.

In this experiment, over 3,500 human participants were paid to provide subjective quality scores, each of whom viewed 53 sequences on average. This ensures more than 160 raw subjective scores collected for each test (non-pristine reference or transcoded) sequence.

B. Data Processing, Validation and Analysis

Due to the nature of crowdsourcing experiments, subjective data acquired in this experiment may be associated with larger variances. Based on the recommendation in [53], we further improved the data reliability by soft screening the collected raw opinion scores when calculating the Mean Opinion Score (MOS), as shown in Algorithm 1. The histograms of the resulting MOS values, the corresponding standard errors (SE), subject inconsistency (σ) and the subject bias (b) are plotted in Fig. 6.

To further validate the reliability of the subjective data collected, based on previous works [59, 70, 77], the raw subjective scores collected for each sequence were randomly divided into two equal groups. For each group, one MOS is calculated for each sequence based on Algorithm 1. We then calculated the Spearman Ranked Order Correlation Coefficient (SROCC) between two groups of MOS for the whole database. This partitioning process has been repeated for 1,000 times to obtain the average SROCC, which is 0.9322. We also provide the scatter plot for one partition in Fig. 7.

Based on the MOS obtained, we can study the influence of reference distortions and different codecs on subjective perceptual experience. We plot the MOS of the transcoded sequences against different quantization parameters for different reference quality levels and transcoding codecs in Fig. 8. Here, we only focus on the HD content encoding without resolution adaptation. It can be observed that when the reference content is associated with relatively high quality (i.e. reference

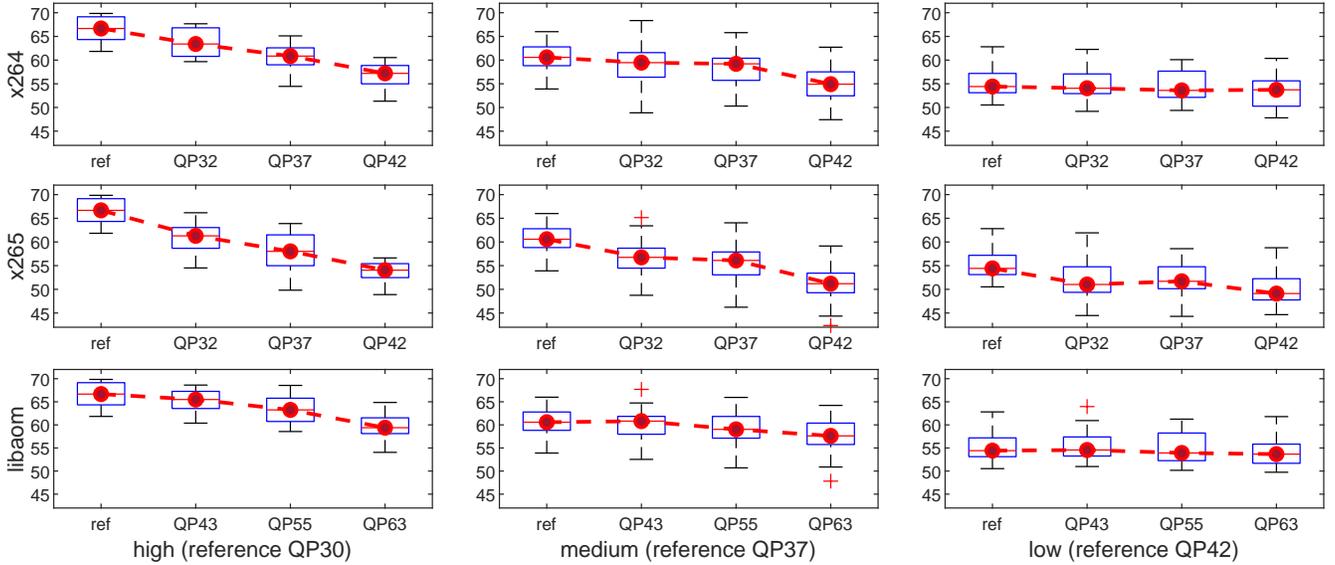


Fig. 8: Boxplots of MOS of the tested sequences against different quantization parameters during transcoding. Three rows correspond to the three transcoding codecs used. Three columns correspond to the three reference quality groups.

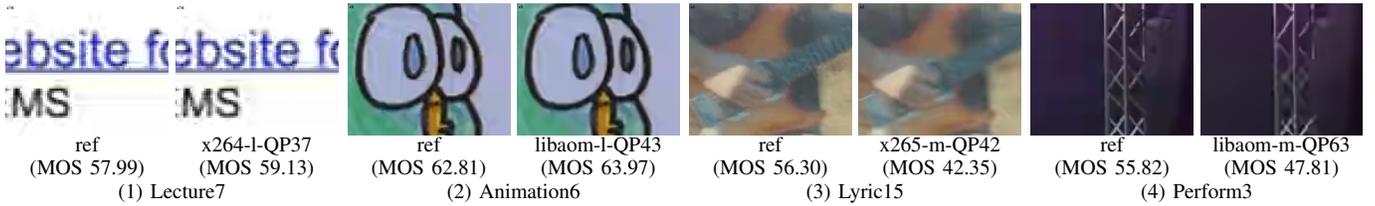


Fig. 9: Visual comparison examples between reference and transcoded content. The labels below figures show the corresponding *codec-group-QP*. For example, x265-m-QP32 means this sample comes from a sequence in the medium reference group, transcoded with x265 codec using QP32.

QP30), applying different QP values during transcoding can generate results with more distinct visual quality levels. In contrast, when the quality of the reference content is low (reference QP42), the visual difference between transcoded content (based on various QP values) is much smaller. We have also provided visual comparison examples between reference and transcoded content in Fig. 9, where in cases (1) and (2), the artifact around the edges were smoothed out during transcoding, which actually improves the visual experience slightly.

V. EVALUATION OF OBJECTIVE QUALITY METRICS

The BVI-UGC database with its associated subjective data contributes a rigorous and valuable benchmarking tool for evaluating the performance of quality assessment methods on user-generated content. More importantly, due to its unique design, it can be used to test both full-reference and no-reference metrics in the context of transcoding. This is achieved by calculating the correlation between quality indices generated by the objective quality models and the DMOS of transcoded content. It can also be employed, like many other UGC databases such as YouTube UGC [3], KoNViD-1k [5, 6], KoNViD-150k [7, 8], and LIVE-VQC [9], to assess the performance of no-reference metrics directly (against MOS) without

references, simulating the quality of experience estimation at the user end.

A. The benchmarked VQA methods

In this section, we present a comprehensive evaluation with 10 popular full-reference (FR) and 11 no-reference (NR) metrics. Specifically, the tested FR VQA methods include PSNR, SSIM [22], MS-SSIM [23], ST-GREED [78], LPIPS [36], VMAF [79], C3DVQA [37], FR-CUGCVQA [38], FR-CONTRIQUE [40], RankDVQA [39] and RankDVQA-UGC [42]. Among these methods, PSNR, SSIM, MS-SSIM are conventional quality metrics, VMAF and ST-GREED are regression-based models, and LPIPS, CONTRIQUE, CUGCVQA, C3DVQA and RankDVQA are deep video quality assessment. It is noted that, among the metrics above, CUGCVQA and CONTRIQUE have both full-reference and no-reference implementations. All FR quality metrics tested here are used to calculate the quality differences between the transcoded sequences and their corresponding non-pristine reference sequences. The Spearman Ranking Order Correlation Coefficients (SROCC) and Kendall Ranking Correlation Coefficients (KRCC) between the predicted quality difference values and their corresponding DMOS are then computed to measure their performance.

TABLE II: The correlation results with DMOS of the tested full-/no-reference metrics on BVI-UGC across different reference and codec groups. For no-reference metrics, the difference between quality indices of transcoded video and its non-pristine reference is taken to correlate with DMOS. In each group, the best and the second best metrics are **boldfaced** and underlined.

Databases	BVI-UGC		BVI-UGC low		BVI-UGC medium		BVI-UGC high		BVI-UGC H264		BVI-UGC H265		BVI-UGC libaom	
	SROCC	KRCC	SROCC	KRCC	SROCC	KRCC	SROCC	KRCC	SROCC	KRCC	SROCC	KRCC	SROCC	KRCC
Full-reference														
PSNR	0.5003	0.3438	0.1913	0.1279	0.5478	0.3795	0.5459	0.3827	0.5335	0.3694	0.4957	0.3409	0.5731	0.4016
SSIM	0.4663	0.3212	0.1978	0.1328	0.5453	0.3766	0.7004	0.5086	0.5235	0.3644	0.5222	0.3629	0.4733	0.3283
ST-GREED [78]	0.1797	0.1307	0.1818	0.1254	0.2237	0.1626	0.1901	0.1325	0.0857	0.0629	0.1197	0.0916	0.3385	0.2388
LPIPS [36]	0.1398	0.0985	0.1731	0.1231	0.1828	0.1299	0.1444	0.0822	0.0564	0.0374	0.0833	0.0580	0.3207	0.2244
VMAF 0.6.1 [79]	<u>0.5610</u>	<u>0.3903</u>	<u>0.4340</u>	<u>0.2980</u>	0.6538	0.4604	<u>0.7588</u>	<u>0.5571</u>	0.6392	0.4521	<u>0.6258</u>	<u>0.4381</u>	<u>0.5948</u>	<u>0.4184</u>
C3DVQA [37]	0.4360	0.2975	0.2678	0.1819	0.4523	0.3035	0.5816	0.4053	0.4810	0.3301	0.4716	0.3194	0.4573	0.3154
FR-CUGCVQA [38]	0.4126	0.2836	0.1837	0.1261	0.4914	0.3334	0.6495	0.4597	0.5456	0.3825	0.4986	0.3475	0.4434	0.3120
FR-CONTRIQUE [40]	0.2150	0.1561	0.1179	0.0889	0.1544	0.1186	0.2686	0.1916	0.6921	0.4924	0.4820	0.3290	0.3725	0.2552
RankDVQA [39]	0.5527	0.3842	0.1405	0.0971	0.5118	0.3514	0.7502	0.5567	0.5920	0.4244	0.6055	0.4283	0.5714	0.4080
RankDVQA-UGC [42]	0.5727	0.4042	0.4420	0.3037	0.6478	0.4574	0.7617	0.5663	0.6347	0.4513	0.6351	0.4403	0.5971	0.4193
No-reference (measuring quality degradation, to correlate with DMOS)														
NIQE [80]	0.1894	0.1247	0.2381	0.1591	0.2589	0.1714	0.1979	0.1301	0.2099	0.1419	0.0236	0.0165	0.2439	0.1642
BRISQUE [81]	0.1316	0.0854	<u>0.2516</u>	<u>0.1628</u>	0.0548	0.0364	0.0824	0.0630	0.1144	0.0742	0.1334	0.0872	0.2210	0.1430
VBLIINDS [82]	0.0651	0.0446	0.0499	0.0367	0.0171	0.0117	0.0542	0.0372	0.2906	0.2067	0.2146	0.1439	0.3217	0.2221
VIIDEO [83]	0.0141	0.0075	0.0918	0.0756	0.0124	0.0126	0.0484	0.0300	0.0349	0.0178	0.0776	0.0547	0.2104	0.1389
ChipQA [50]	0.2418	0.1749	0.1901	0.1342	0.3087	0.2253	0.2894	0.2071	0.2112	0.1532	0.2099	0.1534	0.3780	0.2708
NR-CUGCVQA [38]	<u>0.5063</u>	<u>0.3463</u>	0.0989	0.0638	0.3045	0.1978	<u>0.5790</u>	<u>0.4068</u>	<u>0.6470</u>	<u>0.4567</u>	0.6467	0.4613	<u>0.4927</u>	<u>0.3450</u>
NR-CONTRIQUE [40]	0.3268	0.2208	0.2645	0.1797	0.3390	0.2315	0.3827	0.2582	0.4485	0.3038	0.3458	0.2299	0.3890	0.2617
SimpleVQA [51]	0.5390	0.3741	0.1512	0.1043	0.5218	0.3595	0.5862	0.4200	0.5877	0.4132	0.5494	0.3839	0.6097	0.4221
FastVQA [10]	0.0232	0.0163	0.0443	0.0295	0.0465	0.0301	0.0060	0.0023	0.0065	0.0045	0.0442	0.0316	0.0067	0.0038
FasterVQA [11]	0.0804	0.0535	0.0160	0.0097	0.0741	0.0487	0.0110	0.0060	0.0856	0.0573	0.0752	0.0500	0.1026	0.0685
CONVIQT [52]	0.4934	0.3374	0.1630	0.1107	<u>0.3760</u>	<u>0.2530</u>	0.4857	0.3304	0.7132	0.5266	<u>0.6028</u>	<u>0.4145</u>	0.3735	0.2525

TABLE III: The correlation results with MOS of the tested no-reference metrics on BVI-UGC across different reference and codec groups. In each group, the best and the second best metrics are **boldfaced** and underlined.

Databases	BVI-UGC		BVI-UGC low		BVI-UGC medium		BVI-UGC high		BVI-UGC H264		BVI-UGC H265		BVI-UGC libaom	
	SROCC	KRCC	SROCC	KRCC	SROCC	KRCC	SROCC	KRCC	SROCC	KRCC	SROCC	KRCC	SROCC	KRCC
No-reference (measuring absolute quality, to correlate with MOS)														
NIQE [80]	0.2533	0.1700	0.3816	0.2558	0.2794	0.1871	0.1117	0.0737	0.4077	0.2837	0.3184	0.2155	0.2266	0.1532
BRISQUE [81]	0.2527	0.1758	0.2313	0.1465	0.3215	0.2200	0.3333	0.2318	0.2498	0.1761	0.2356	0.1667	0.2860	0.1991
VBLIINDS [82]	0.1211	0.0808	0.0565	0.0361	0.1355	0.0911	0.0143	0.0106	0.1384	0.0956	0.1516	0.1017	0.3513	0.2370
VIIDEO [83]	0.1651	0.1102	0.2878	0.1940	0.1322	0.0881	0.0261	0.0181	0.1565	0.1055	0.2904	0.1943	0.1450	0.0951
ChipQA [50]	0.1994	0.1464	0.0652	0.0419	0.2687	0.1951	0.3143	0.2313	0.1754	0.1292	0.1642	0.1212	0.3040	0.2245
NR-CUGCVQA [38]	0.4333	0.2953	0.0664	0.0469	0.3649	0.2479	<u>0.4723</u>	<u>0.3191</u>	0.6202	0.4433	<u>0.5548</u>	<u>0.3904</u>	0.2940	0.1977
NR-CONTRIQUE [40]	0.3376	0.2294	0.2076	0.1417	0.3738	0.2552	0.3056	0.2113	0.3023	0.2045	0.3514	0.2420	<u>0.4321</u>	<u>0.2977</u>
SimpleVQA [51]	0.5203	0.3641	0.1439	0.0983	0.5562	0.3933	0.5345	0.3702	<u>0.5832</u>	<u>0.4186</u>	0.5391	0.3821	0.5586	0.3926
FastVQA [10]	0.1897	0.1266	0.3417	0.2236	0.2905	0.1951	0.1539	0.1033	0.2186	0.1455	0.2091	0.1405	0.1974	0.1333
FasterVQA [11]	0.1922	0.1285	0.3215	0.2051	0.2933	0.1966	0.1470	0.0984	0.2223	0.1483	0.2116	0.1422	0.1889	0.1285
CONVIQT [52]	<u>0.4678</u>	<u>0.3232</u>	<u>0.3507</u>	<u>0.2323</u>	<u>0.4855</u>	<u>0.3359</u>	0.3594	0.2469	0.5387	0.3773	0.5590	0.3927	0.4002	0.2750

Similarly, 11 NR quality metrics also include conventional or machine learning-based methods such as NIQE [80], BRISQUE [81], VBLIINDS [82], VIIDEO [83], and deep learning-based quality models, e.g., ChipQA [50], NR-CUGCVQA [38], NR-CONTRIQUE [40], SimpleVQA [51], CONVIQT [52], FastVQA [10] and FasterVQA [11]. For each NR metric, we first adapt to the FR transcoding scenario by calculating the quality indices of transcoded videos and its non-pristine reference separately, and obtaining their quality differences. This is to measure the quality degradation of the transcoded content relative to the reference. We also assess the performance of these NR models by calculating their correlation coefficients between their predicted quality indices and the corresponding MOS of the transcoded sequences.

For all the learning-based methods, including LPIPS, VMAF, ChipQA, CUGCVQA, CONTRIQUE, ST-GREED, C3DVQA, RankDVQA, BRISQUE, SimpleVQA, CONVIQT, FastVQA and FasterVQA, their pre-trained models provided in the associated original literature were used in this experiment. To test model generalization, we did not perform any cross-validation within the proposed database.

B. Evaluation results

TABLE II summarizes the performance results of all the tested full-/no-reference VQA methods in the context of transcoding, where the DMOS values are employed to calculate the correlation coefficients. It is noted that none of the tested quality assessment methods achieve satisfactory overall correlation performance on BVI-UGC - the best performer RankDVQA-UGC only offers a SROCC value of 0.5727. Among all the no-reference quality metrics, SimpleVQA achieves the highest SROCC value both when measuring the quality degradation (0.5390) and when predicting the absolute quality (0.5203). We have further divided the whole BVI-UGC database into three subsets according to the unpristine references. It can be observed that for test sequences generated from low-quality references, most full-reference quality models perform worse than on those from medium-/high-quality references. This confirms our assumption that the quality of reference videos does affect the quality prediction accuracy of full-reference quality models. We performed another segmentation of the database based on the codec used in transcoding, and found that many quality metrics achieve better performance on H.264 or H.265 content compared to libaom compressed sequences. This may be because of the

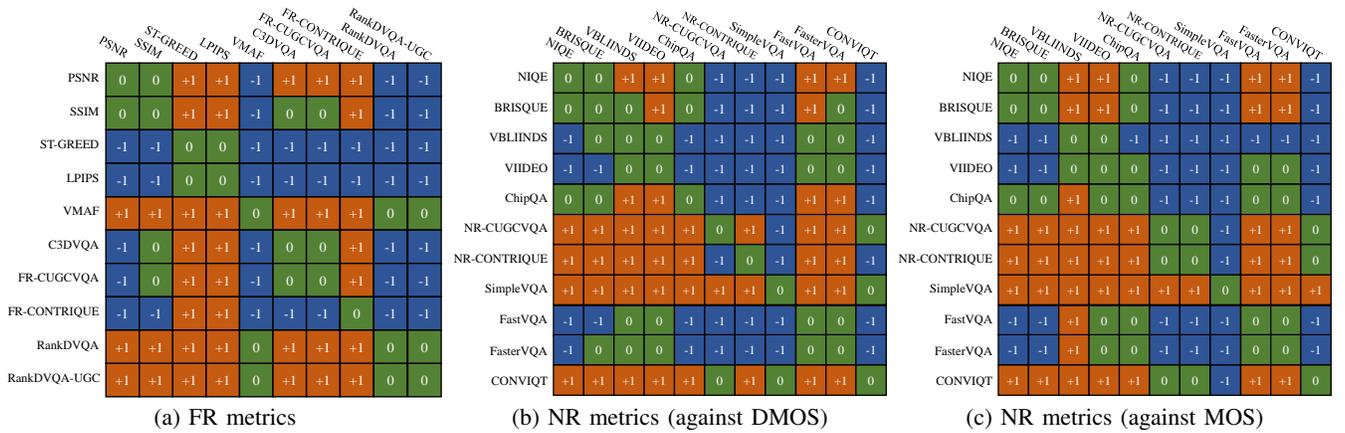


Fig. 10: Pairwise comparisons between the overall performances. The color and value of the cells indicate the F-test results between the DMOS prediction residuals of the metrics pair, at a 95% confidence interval. Orange cell with +1 value indicates the metric in the row is superior to the metric in the column and blue cell with -1 value means the opposite. Green cell with 0 value denotes statistical equivalence.

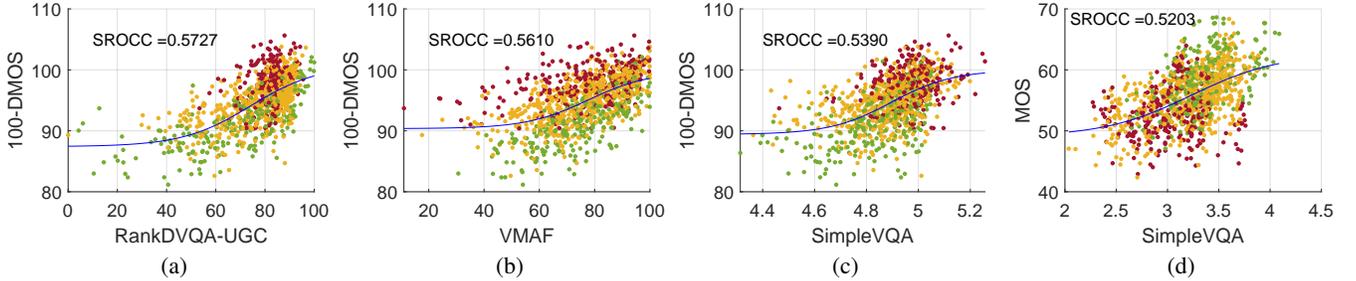


Fig. 11: Scatter plots of subjective scores against the predictions of the best-performing models. The green, yellow, and red scatter points correspond to the transcoded sequences generated from three groups with different reference qualities, high, medium, and low. The blue lines are the logistic functions fitted between the model predictions and subjective scores on the entire BVI-UGC database.

fact that these VQA methods are trained and/or validated more often on H.264/H.265 compressed content.

As mentioned above, the BVI-UGC database can also be used to evaluate no-reference quality metrics in terms of their ability to directly predict the visual quality of distorted videos (using MOS to calculate correlation coefficients). TABLE III summarizes their results when benchmarked on the BVI-UGC database. Here the best overall performance is also provided by SimpleVQA, while the second best performer is CONVIQT. Based on all these results we can conclude that assessing the perceptual quality of UGC content is a highly challenging task, in particular when the reference content is distorted. More advanced VQA models are urgently required to offer enhanced prediction performance.

To further validate the performance ranking among the benchmarked quality assessment methods, we also conducted an F-test between every two metrics to check the statistical significance of their difference, following the practice in [30, 55]. Specifically, a pairwise comparison was performed on the residuals between the DMOS (or MOS) and the model prediction after non-linear regression [84]. Fig. 10 (a-c) summarize the F-test results between full-/no-reference metrics (on DMOS) and no-reference models (on MOS). We can observe that RankDVQA-UGC, RankDVQA and VMAF are the best performers among full-reference metrics, all of which significantly outperform seven other quality models in a 95% confidence interval. For no-reference quality models,

SimpleVQA is statistically better than ten other NR VQA methods based on DMOS or MOS.

In order to analyze the correlation performance of the objective quality metrics, the scatter plots of the predictions of the selected, well-performing models against subjective scores, along with the fitted logistic curves, are shown in Fig. 11. It can be observed that, in all cases, the scatter points are distributed sparsely along the fitting curves, which also demonstrates the unsatisfactory performance of existing VQA methods on this database from a different perspective.

All the results shown in this section confirm the urgent need for an accurate and robust quality metric, which can adapt to various reference content scenarios and different distortion types, to facilitate UGC streaming applications.

VI. CONCLUSION

In this paper, we have presented a novel video database, BVI-UGC, which is the first UGC database to contain (non-pristine) references with various levels of distortions and transcoded content generated by multiple codecs. It consists of 60 pseudo-pristine source sequences with diverse and representative user-generated video content, covering 15 popular UGC categories. These were further used to produce 60 non-pristine reference sequences and 1,080 transcoded sequences following a typical UGC streaming pipeline. Based on this database, we designed and performed a large-scale crowd-sourcing subjective study on the perceptual quality of of both

non-pristine referece and transcoded videos. The collected subjective scores, together with the video clips, have been employed to benchmark the performance of 21 full-reference and no-reference popular quality assessment methods. The results clearly show that all these quality metrics fail to perform well on this database, with ranking order correlation coefficient (SROCC) values below 0.6.

We believe that the BVI-UGC database will provide a valuable resource to the research community for developing and validating new video quality assessment models in the context of UGC transcoding. Future work is now required to investigate full-reference and no-reference quality metrics, which can predict the perceived quality of streamed UGC content more accurately and robustly.

REFERENCES

- [1] Sandvine, “Global internet phenomena report, 2024,” 2024. [Online]. Available: <https://www.sandvine.com/>
- [2] L. Merritt and R. Vanam, “x264: A high performance h.264/avc encoder,” 2006. [Online]. Available: http://neuron2.net/library/avc/overview_x264_v8_5.pdf
- [3] Y. Wang, S. Inguva, and B. Adsumilli, “Youtube UGC dataset for video compression research,” in *2019 IEEE 21st International Workshop on Multimedia Signal Processing (MMSP)*. IEEE, 2019, pp. 1–5.
- [4] Z. Tu, Y. Wang, N. Birkbeck, B. Adsumilli, and A. C. Bovik, “UGC-VQA: Benchmarking blind video quality assessment for user generated content,” *IEEE Trans. on Image Processing*, vol. 30, pp. 4449–4464, 2021.
- [5] V. Hosu, F. Hahn, M. Jenadeleh, H. Lin, H. Men, T. Szirányi, S. Li, and D. Saupe, “The konstanz natural video database,” 2017. [Online]. Available: <https://database.mmsp-kn.de>
- [6] V. Hosu, F. Hahn, M. Jenadeleh, H. Lin, H. Men, T. Szirányi, S. Li, and D. Saupe, “The konstanz natural video database (konvid-1k),” in *2017 Ninth International Conf. on Quality of Multimedia Experience (QoMEX)*. IEEE, 2017, pp. 1–6.
- [7] F. Götz-Hahn, V. Hosu, H. Lin, and D. Saupe, “The konstanz 150k in-the-wild video database (konvid-150k),” 2021. [Online]. Available: <https://database.mmsp-kn.de>
- [8] F. Götz-Hahn, V. Hosu, H. Lin, and D. Saupe, “Konvid-150k: A dataset for no-reference video quality assessment of videos in-the-wild,” in *IEEE Access* 9. IEEE, 2021, pp. 72 139–72 160.
- [9] Z. Sinno and A. C. Bovik, “Large-scale study of perceptual video quality,” *IEEE Trans. on Image Processing*, vol. 28, no. 2, pp. 612–627, 2018.
- [10] H. Wu, C. Chen, J. Hou, L. Liao, A. Wang, W. Sun, Q. Yan, and W. Lin, “Fast-VQA: Efficient end-to-end video quality assessment with fragment sampling,” in *European conference on computer vision*. Springer, 2022, pp. 538–554.
- [11] H. Wu, C. Chen, L. Liao, J. Hou, W. Sun, Q. Yan, J. Gu, and W. Lin, “Neighbourhood representative sampling for efficient end-to-end video quality assessment,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [12] Y. Li, S. Meng, X. Zhang, S. Wang, Y. Wang, and S. Ma, “UGC-VIDEO: perceptual quality assessment of user-generated videos,” in *2020 IEEE Conf. on Multimedia Information Processing and Retrieval (MIPR)*. IEEE, 2020, pp. 35–38.
- [13] X. Yu, N. Birkbeck, Y. Wang, C. G. Bampis, B. Adsumilli, and A. C. Bovik, “Predicting the quality of compressed videos with pre-existing distortions,” *IEEE Trans. on Image Processing*, vol. 30, pp. 7511–7526, 2021.
- [14] H. Wang, G. Li, S. Liu, and C.-C. J. Kuo, “Challenge on quality assessment of compressed UGC videos,” 2021. [Online]. Available: https://2021.ieeeicme.org/2021.ieeeicme.org/conf_challenges.html
- [15] “x265, code repository.” [Online]. Available: <https://github.com/videoan/x265>
- [16] “Alliance for Open Media.” [Online]. Available: <https://aomedia.org/>
- [17] A. M. Turk, “Amazon mechanical turk,” *Retrieved Dec*, 2023.
- [18] D. Bull and F. Zhang, *Intelligent image and video compression: communicating pictures*. Academic Press, 2021.
- [19] J. A. Ferwerda, P. Shirley, S. N. Pattanaik, and D. P. Greenberg, “A model of visual masking for computer graphics,” in *Proceedings of the 24th annual conference on Computer graphics and interactive techniques*, 1997, pp. 143–152.
- [20] X. Yang, W. Ling, Z. Lu, E. P. Ong, and S. Yao, “Just noticeable distortion model and its applications in video coding,” *Signal processing: Image communication*, vol. 20, no. 7, pp. 662–680, 2005.
- [21] A. P. Ginsburg, “Contrast sensitivity and functional vision,” *International ophthalmology clinics*, vol. 43, no. 2, pp. 5–15, 2003.
- [22] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, “Image quality assessment: from error visibility to structural similarity,” *IEEE Trans. on image processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [23] Z. Wang, E. P. Simoncelli, and A. C. Bovik, “Multiscale structural similarity for image quality assessment,” in *The Thirty-Seventh Asilomar Conf. on Signals, Systems & Computers*, 2003, vol. 2. Ieee, 2003, pp. 1398–1402.
- [24] A. K. Moorthy and A. C. Bovik, “A motion compensated approach to video quality assessment,” in *2009 Conference Record of the Forty-Third Asilomar Conference on Signals, Systems and Computers*. IEEE, 2009, pp. 872–875.
- [25] A. K. Moorthy and A. C. Bovik, “Efficient video quality assessment along temporal trajectories,” *IEEE transactions on circuits and systems for video technology*, vol. 20, no. 11, pp. 1653–1658, 2010.
- [26] K. Zeng and Z. Wang, “3D-SSIM for video quality assessment,” in *2012 19th IEEE international conference on image processing*. IEEE, 2012, pp. 621–624.
- [27] A. Rehman, K. Zeng, and Z. Wang, “Display device-adapted video quality-of-experience assessment,” in *Human vision and electronic imaging XX*, vol. 9394. SPIE, 2015, pp. 27–37.
- [28] S. Li, F. Zhang, L. Ma, and K. N. Ngan, “Image quality assessment by separately evaluating detail losses and additive impairments,” *IEEE Trans. on Multimedia*, vol. 13, no. 5, pp. 935–949, 2011.
- [29] H. R. Sheikh and A. C. Bovik, “Image information and visual quality,” *IEEE Transactions on image processing*, vol. 15, no. 2, pp. 430–444, 2006.
- [30] K. Seshadrinathan and A. C. Bovik, “Motion tuned spatio-temporal quality assessment of natural videos,” *IEEE Trans. on image processing*, vol. 19, no. 2, pp. 335–350, 2009.
- [31] P. V. Vu, C. T. Vu, and D. M. Chandler, “A spatiotemporal most-apparent-distortion model for video quality assessment,” in *2011 18th IEEE International Conf. on Image Processing*. IEEE, 2011, pp. 2505–2508.
- [32] F. Zhang and D. R. Bull, “A perception-based hybrid model for video quality assessment,” *IEEE Trans. on Circuits and Systems for Video Technology*, vol. 26, no. 6, pp. 1017–1028, 2015.
- [33] C. Cortes and V. Vapnik, “Support-vector networks,” *Machine learning*, vol. 20, no. 3, pp. 273–297, 1995.
- [34] F. Zhang, A. Katsenou, C. Bampis, L. Krasula, Z. Li, and D. Bull, “Enhancing vmaf through new feature integration and model combination,” in *2021 Picture Coding Symposium (PCS)*. IEEE, 2021, pp. 1–5.
- [35] W. Kim, J. Kim, S. Ahn, J. Kim, and S. Lee, “Deep video quality assessor: From spatio-temporal visual sensitivity to a convolutional neural aggregation network,” in *Proc. of the European Conf. on Computer Vision*, 2018, pp. 219–234.
- [36] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, “The unreasonable effectiveness of deep features as a perceptual

- metric,” in *Proc. of the IEEE Conf. on computer vision and pattern recognition*, 2018, pp. 586–595.
- [37] M. Xu, J. Chen, H. Wang, S. Liu, G. Li, and Z. Bai, “C3DVQA: Full-reference video quality assessment with 3d convolutional neural network,” in *ICASSP 2020-2020 IEEE International Conf. on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 4447–4451.
- [38] Y. Li, L. Feng, J. Xu, T. Zhang, Y. Liao, and J. Li, “Full-reference and no-reference quality assessment for compressed user-generated content videos,” in *2021 IEEE International Conf. on Multimedia & Expo Workshops (ICMEW)*. IEEE, 2021, pp. 1–6.
- [39] C. Feng, D. Danier, F. Zhang, and D. Bull, “RankDVQA: Deep vqa based on ranking-inspired hybrid training,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024, pp. 1648–1658.
- [40] P. C. Madhusudana, N. Birkbeck, Y. Wang, B. Adsumilli, and A. C. Bovik, “Image quality assessment using contrastive learning,” *IEEE Transactions on Image Processing*, vol. 31, pp. 4149–4161, 2022.
- [41] T. Peng, C. Feng, D. Danier, F. Zhang, and D. Bull, “RMT-BVQA: Recurrent memory transformer-based blind video quality assessment for enhanced video content,” *arXiv preprint arXiv:2405.08621*, 2024.
- [42] Z. Qi, C. Feng, D. Danier, F. Zhang, X. Xu, S. Liu, and D. Bull, “Full-reference video quality assessment for user generated content transcoding,” *arXiv preprint arXiv:2312.12317*, 2023.
- [43] K. Zhu, C. Li, V. Asari, and D. Saupe, “No-reference video quality assessment based on artifact measurement and statistical analysis,” *IEEE Trans. on Circuits and Systems for Video Technology*, vol. 25, no. 4, pp. 533–546, 2014.
- [44] J. E. Caviedes and F. Oberti, “No-reference quality metric for degraded and enhanced video,” in *Digital Video Image Quality and Perceptual Coding*. CRC Press, 2017, pp. 305–324.
- [45] F. Zhang, W. Lin, Z. Chen, and K. N. Ngan, “Additive log-logistic model for networked video quality assessment,” *IEEE Trans. on Image Processing*, vol. 22, no. 4, pp. 1536–1547, 2012.
- [46] Z. Wang, W. Wang, Z. Wan, Y. Xia, and W. Lin, “No-reference hybrid video quality assessment based on partial least squares regression,” *Multimedia tools and applications*, vol. 74, pp. 10277–10290, 2015.
- [47] Y. Wang, S.-U. Kum, C. Chen, and A. Kokaram, “A perceptual visibility metric for banding artifacts,” in *2016 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2016, pp. 2067–2071.
- [48] D. Ghadiyaram, C. Chen, S. Inguva, and A. Kokaram, “A no-reference video quality predictor for compression and scaling artifacts,” in *2017 IEEE International Conf. on Image Processing (ICIP)*. IEEE, 2017, pp. 3445–3449.
- [49] W. Liu, Z. Duanmu, and Z. Wang, “End-to-end blind quality assessment of compressed videos using deep neural networks,” in *ACM Multimedia*, 2018, pp. 546–554.
- [50] J. P. Ebenezer, Z. Shang, Y. Wu, H. Wei, S. Sethuraman, and A. C. Bovik, “ChipQA: No-reference video quality prediction via space-time chips,” *IEEE Transactions on Image Processing*, vol. 30, pp. 8059–8074, 2021.
- [51] W. Sun, X. Min, W. Lu, and G. Zhai, “A deep learning based no-reference quality assessment model for UGC videos,” in *Proc. of the 30th ACM International Conf. on Multimedia*, 2022, pp. 856–865.
- [52] P. C. Madhusudana, N. Birkbeck, Y. Wang, B. Adsumilli, and A. C. Bovik, “Conviqt: Contrastive video quality estimator,” *IEEE Transactions on Image Processing*, 2023.
- [53] I. Recommendation, “910,” subjective video quality assessment methods for multimedia applications,” recommendation ITU-T P. 910,” *ITU Telecom. Standardization Sector of ITU*, 2022.
- [54] “VQEG: The video quality experts group.” [Online]. Available: <https://www.vqeg.org>
- [55] K. Seshadrinathan, R. Soundararajan, A. C. Bovik, and L. K. Cormack, “Study of subjective and objective quality assessment of video,” *IEEE Trans. on Image Processing*, vol. 19, no. 6, pp. 1427–1441, 2010.
- [56] F. Zhang, S. Li, L. Ma, Y. C. Wong, and K. N. Ngan, “IVP subjective quality video database,” *The Chinese University of Hong Kong*, <https://ivp.ee.cuhk.edu.hk/research/database/subjective>, 2011.
- [57] F. Zhang, F. M. Moss, R. Baddeley, and D. R. Bull, “BVI-HD: A video quality database for HEVC compressed and texture synthesized content,” *IEEE Trans. on Multimedia*, vol. 20, no. 10, pp. 2620–2630, 2018.
- [58] A. Mackin, F. Zhang, and D. R. Bull, “A study of subjective video quality at various frame rates,” in *2015 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2015, pp. 3407–3411.
- [59] P. C. Madhusudana, X. Yu, N. Birkbeck, Y. Wang, B. Adsumilli, and A. C. Bovik, “Subjective and objective quality assessment of high frame rate videos,” *IEEE Access*, vol. 9, pp. 108 069–108 082, 2021.
- [60] P. C. Madhusudana, N. Birkbeck, Y. Wang, B. Adsumilli, and A. C. Bovik, “Capturing video frame rate variations through entropic differencing,” *arXiv e-prints*, pp. arXiv–2006, 2020.
- [61] A. Mackin, M. Afonso, F. Zhang, and D. Bull, “A study of subjective video quality at various spatial resolutions,” in *2018 25th IEEE International Conference on Image Processing (ICIP)*. IEEE, 2018, pp. 2830–2834.
- [62] A. Mackin, D. Ma, F. Zhang, and D. Bull, “A subjective study on videos at various bit depths,” *arXiv preprint arXiv:2103.10363*, 2021.
- [63] Z. Chen, W. Sun, Z. Zhang, R. Huang, F. Lu, X. Min, G. Zhai, and W. Zhang, “FS-BAND: A frequency-sensitive banding detector,” *arXiv preprint arXiv:2311.18216*, 2023.
- [64] Z. Chen, W. Sun, J. Jia, F. Lu, Z. Zhang, J. Liu, R. Huang, X. Min, and G. Zhai, “BAND-2k: Banding artifact noticeable database for banding detection and quality assessment,” *IEEE Transactions on Circuits and Systems for Video Technology*, pp. 1–16, 2024.
- [65] Z. Zhang, W. Wu, W. Sun, D. Tu, W. Lu, X. Min, Y. Chen, and G. Zhai, “MD-VQA: Multi-dimensional quality assessment for UGC live videos,” in *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2023, pp. 1746–1755.
- [66] F. M. Moss, K. Wang, F. Zhang, R. Baddeley, and D. R. Bull, “On the optimal presentation duration for subjective video quality assessment,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 26, no. 11, pp. 1977–1987, 2015.
- [67] F. M. Moss, F. Zhang, R. Baddeley, and D. R. Bull, “What’s on TV: A large scale quantitative characterisation of modern broadcast video content,” in *2016 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2016, pp. 2425–2429.
- [68] D. Ma, F. Zhang, and D. Bull, “BVI-DVC: A training database for deep video compression,” *IEEE Trans. on Multimedia*, 2021.
- [69] S. Winkler, “Analysis of public image and video databases for quality assessment,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 6, no. 6, pp. 616–625, 2012.
- [70] D. Danier, F. Zhang, and D. R. Bull, “BVI-VFI: A video quality database for video frame interpolation,” *IEEE Transactions on Image Processing*, 2023.
- [71] I. Telecom, “Advanced video coding for generic audiovisual services,” *ITU-T Recommendation H. 264*, 2003.
- [72] “FFmpeg.” [Online]. Available: <https://ffmpeg.org/>
- [73] (2020) AOM video model (AVM). [Online]. Available: <https://gitlab.com/AOMediaCodec/avm>
- [74] X. Zhao, Z. Lei, A. Norkin, T. Daede, and A. Tourapis, “AOM common test conditions v2. 0,” *Alliance for Open Media, Codec Working Group Output Document*, 2021.
- [75] R. I.-R. BT, “Methodology for the subjective assessment of the quality of television pictures,” *International Telecommunication Union*, 2002.

- [76] “Infobyip.com.” [Online]. Available: <https://www.infobyip.com/detectdisplaysize.php>
- [77] T. Hoßfeld, C. Keimel, M. Hirth, B. Gardlo, J. Habigt, K. Diepold, and P. Tran-Gia, “Best practices for QoE crowdtesting: QoE assessment with crowdsourcing,” *IEEE transactions on multimedia*, vol. 16, no. 2, pp. 541–558, 2013.
- [78] P. C. Madhusudana, N. Birkbeck, Y. Wang, B. Adsumilli, and A. C. Bovik, “ST-GREED: Space-time generalized entropic differences for frame rate dependent video quality prediction,” *IEEE Transactions on Image Processing*, vol. 30, pp. 7446–7457, 2021.
- [79] Z. Li, A. Aaron, I. Katsavounidis, A. Moorthy, and M. Manohara, “Toward a practical perceptual video quality metric,” *The Netflix Tech Blog*, vol. 6, no. 2, 2016.
- [80] A. Mittal, R. Soundararajan, and A. C. Bovik, “Making a “completely blind” image quality analyzer,” *IEEE Signal processing letters*, vol. 20, no. 3, pp. 209–212, 2012.
- [81] A. Mittal, A. K. Moorthy, and A. C. Bovik, “No-reference image quality assessment in the spatial domain,” *IEEE Trans. on image processing*, vol. 21, no. 12, pp. 4695–4708, 2012.
- [82] M. A. Saad, A. C. Bovik, and C. Charrier, “Blind prediction of natural video quality,” *IEEE Trans. on image Processing*, vol. 23, no. 3, pp. 1352–1365, 2014.
- [83] A. Mittal, M. A. Saad, and A. C. Bovik, “A completely blind video integrity oracle,” *IEEE Trans. on Image Processing*, vol. 25, no. 1, pp. 289–300, 2015.
- [84] “Final report from the video quality experts group on the validation of objective quality metrics for video quality assessment,” https://www.its.bldrdoc.gov/vqeg/projects/frtv_phase1, 2000.