# DDH-QA: A DYNAMIC DIGITAL HUMANS QUALITY ASSESSMENT DATABASE

*Zicheng Zhang[1,2], Yingjie Zhou[1,2], Wei Sun[1,2], Wei Lu[1,2], Xiongkuo Min[1,2], Yu Wang[1], and Guangtao Zhai[1,2,3]*

[1]Institute of Image Communication and Network Engineering, Shanghai Jiao Tong University, China
[2] Peng Cheng Laboratory, China
[3] MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University, China

## ABSTRACT

In recent years, large amounts of effort have been put into pushing forward the real-world application of dynamic digital human (DDH). However, most current quality assessment research focuses on evaluating static 3D models and usually ignores motion distortions. Therefore, in this paper, we construct a large-scale dynamic digital human quality assessment (DDH-QA) database with diverse motion content as well as multiple distortions to comprehensively study the perceptual quality of DDHs. Both model-based distortion (noise, compression) and motion-based distortion (binding error, motion unnaturalness) are taken into consideration. Ten types of common motion are employed to drive the DDHs and a total of 800 DDHs are generated in the end. Afterward, we render the video sequences of the distorted DDHs as the evaluation media and carry out a well-controlled subjective experiment. Then a benchmark experiment is conducted with the state-of-the-art video quality assessment (VQA) methods and the experimental results show that existing VQA methods are limited in assessing the perceptual loss of DDHs. The database is available at https://github.com/zzc-1998/DDH-QA.

***Index Terms***— Dynamic digital human, model-based distortion, motion-based distortion, subjective experiment

## I. INTRODUCTION

Digital humans indicate digital models represented by computer graphics, which are usually static and fixed. Dynamic digital humans (DDHs) are digital models driven by predefined animations [1], which have been widely adopted in applications such as the game industry, film post-production, metaverse, etc. As shown in Fig. 1, the DDHs suffer from both model-based and motion-based distortions. The model-based distortions represent the degradations directly affecting the digital human models. For example, the slight geometry shift and inevitable camera noise can introduce noise disturbance to the geometry structure and texture maps during the generation procedure [2], [3]. Moreover, to
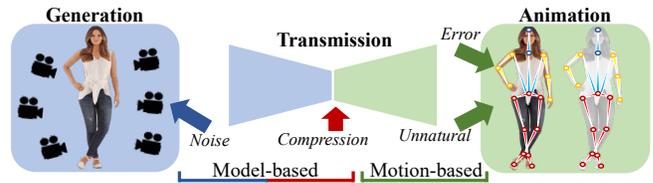
**Fig. 1**. Distortion sources of DDHs.

support real-time VR/AR applications under restricted bandwidth, the digital human models always undergo compression through the transmission procedure. The motion-based distortions stand for the incoherence and unnaturalness of the animation, which are often caused by inappropriate skeleton binding and confusing motion. Nowadays, researchers are mainly paying attention to the generation, representation, rendering, and animation of digital humans [4]. However, the quality assessment for DDHs has fallen behind and effective approaches along with databases are urgently needed.

Therefore, we propose the first dynamic digital human quality assessment (DDH-QA) to tackle the mentioned challenge. One male and one female digital humans represented by texture meshes are collected as the reference. Then we introduce six model-based distortions (color noise, geometry noise, texture compression, texture downsampling, position compression, and UV map compression) and two motion-based distortions (skeleton binding error and motion range unnaturalness) to reference models, which generate 800 distorted DDHs in total. Afterward, we carry out a subjective experiment to collect the perceived quality ratings for the distorted DDHs. Finally, we conduct a benchmark experiment on the DDH-QA database with state-of-the-art video quality assessment (VQA) methods, which shows that current quality assessment methods are not effective for predicting the visual quality levels of DDHs. Our contributions can be summarized as followed:

- **To the best of our knowledge, we construct the first dynamic digital human quality assessment database**, which provides 800 distorted DDHs with both model-based and motion-based distortions.
- We carry out a subjective study to collect the perceptual

quality labels of the distorted DDHs. A total of 32,800 = 41×800 quality ratings are gathered.
- We conduct a benchmark experiment to exhibit the performance of the existing state-of-the-art quality assessment methods.

**Table I**. The comparison of 3D-QA databases and our database, where 'Num' represents the number of the models provided with quality labels.

| Database | Num | Content |
|----------|-----|---------|
| SJTU-PCQA [5] | 378 | Colored Point Cloud |
| WPC [6] | 740 | Colored Point Cloud |
| LSPCQA [7] | 1,240 | Colored Point Cloud |
| CMDM [8] | 80 | Colored Mesh |
| TMQA [9] | 3,000 | Textured Mesh |
| VVDB2 [10] | 152 | Volumetric Video |
| DHHQA [11] | 1,540 | Static Digital Human Head |
| **DDH-QA**(Ours) | 800 | Dynamic Digital Human |

## II. RELATED WORKS

### II-A. 3D Model Quality Assessment Database

In this section, we give a brief review of the 3D model quality assessment (3D-QA) databases. Mainstream 3D-QA databases focus on static point cloud quality assessment (PCQA) and mesh quality assessment (MQA) [12], [13], [14]. Namely, the SJTU-PCQA [5], WPC [6], and LSPCQA [7] databases contain 378, 740, and 1,240 subjective annotated colored point clouds respectively, which are distorted by noise, downsampling, and compression. Some researchers are also interested in MQA tasks. For example, the CMDM [8] database employs simplification and quantization algorithms to obtain 80 distorted colored meshes. The TMQ database further provides 3,000 distorted textured meshes by compressing both the geometry structure and texture maps. All the databases mentioned above are constructed for common 3D objects, then some databases are proposed to focus on 3D digital humans. The VVDB2 [10] database provides 152 3D human volumetric videos and the DHHQA [11] database includes 1,540 distorted static digital human heads. However, none of these databases specifically investigate the perceptual quality of DDHs and all of them ignore the motion distortion.

### II-B. VQA Development

Since most DDHs are presented in the format of rendered 2D videos, it is reasonable to transfer the VQA models to the DDH-QA tasks. The VQA methods can be generally categorized into full-reference (FR) and no-reference (NR) VQA methods according to the availability of the reference videos. The FR-VQA methods usually compare the frame-level difference with the assistance of image quality assessment (IQA) metrics such as PSNR and SSIM [15]. For the NR-VQA development, some handcrafted-based methods [16], [17], [18], [19] are proposed to extract features based
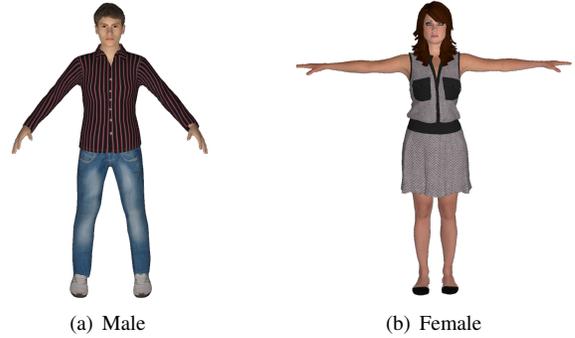


(a) Male              (b) Female

**Fig. 2**. Illustration of the source male and female digital human models. The male digital human model is displayed in 'A' pose while the female digital human model is displayed in 'T' pose respectively.

on natural scene statistics (NSS) and regress the features via the Support Vector Machine. With the development of deep neural networks (DNN), some researchers [20], [21], [22] propose to utilize DNNs for feature extraction and have greatly boosted the performance of NR-VQA models.

## III. DATABASE CONSTRUCTION

### III-A. Source Model Collection

To build the dynamic digital human quality assessment (DDH-QA) database, we collect one male and one female digital human models which can be freely downloaded from **Cgtrader**[1] and are represented by textured meshes. The male model contains 19,528 vertices and 2K texture maps while the female model contains 16,351 vertices and 2K texture maps. The front projections of the digital human models are illustrated in Fig. 2.

### III-B. Distortion Generation

To fit the practical situation of producing DDHs, we introduce two types of distortions, including model-based and motion-based distortions. The overview of the introduced distortion is shown in Table II.

**(1) Model-based Distortion**: The model-based distortions focus on the noise and compression artifacts during the generation and transmission procedure, which include: (a) Color Noise (CN): Gaussian noise is added to the texture maps with $\sigma_c$ set as $\{20, 40, 60, 80, 100\}$; (b) Geometry Noise (GN): Gaussian noise is introduced to the vertices with $\sigma_g$ set as $\{0.01, 0.02, 0.03, 0.04, 0.05\}$; (c) Texture Compression (TC): The texture maps are compressed with JPEG and the quality levels are set as $\{3, 7, 15, 20, 25\}$; (d) Texture Downsampling (TD): The texture maps are downsampled with the sampling rate of $\{2, 4, 8, 12, 16\}$; (e) Position Compression (PC): The position attributes are quantified wit the Draco [23] library and quantization parameters **qp** are set

[1]https://www.cgtrader.com/

**Table II**. Overview of the introduced distortion.

| Type | Distortion | Description |
|------|------------|-------------|
| Model-based | CN | Color noise on the texture maps |
| | GN | Geometry noise on the vertices |
| | TC | Texture maps JPEG compression |
| | TD | Texture maps downsampling |
| | PC | Position quantization by Draco |
| | UMC | Texture coordinate quantization by Draco |
| Motion-based | SBE | Skeleton binding error |
| | MRU | Unnatural motion range |

as {6, 7, 8, 9, 10}; (f) UV Map Compression (UMC): The texture coordinate attributes are quantified with the Draco library and the quantization parameters **qt** are set as {3, 4, 5, 6, 7}; The distorted samples are exhibited in Fig. 3

**(2) Motion-based Distortion:** The motion-based distortions focus on the skeleton rigging bias along with the motion unnaturalness. In most piratical situations, digital humans are first processed with skeleton binding and then animated with the designed motion. The skeleton binding error and the poor-designed motion can cause confusion and unnaturalness to the dynamic digital human animation. Therefore, we introduce the motion-based distortions from two aspects: (a) Skeleton Binding Error (SBE) : Mismatch of skeleton key points are added under five manual defined levels to cover most quality ranges (slight mismatch ~ severe mismatch); (b) Motion Range Unnaturalness (MRU) : We manually adjust the motion range to model the motion unnaturalness under five strengthens. The examples of the motion-based distortions are shown in Fig. 4.

### III-C. Video Rendering

Since the DDHs are usually perceived in the format of 2D animation videos, we decide to render the DDHs into videos for evaluation. We bind the skeleton of the digital human models and render the animation videos with a resolution of 1080P by using Maya software [2] (the viewpoints are manually selected to cover sufficient quality content). To enrich the motion content diversity, we select ten types of common motion, including baseball, boxing, dance, golf, jog, jump, pushup, roll, walk, and wave. The overview of the ten kinds of motion is exhibited in Fig. 5. To sum up, a total of 800 = 2×8×5×10 (digital human models×distortion types×distortion levels×motion types) DDH video sequences are generated for evaluation.

### III-D. Subjective Experiment

We carry out the subjective quality assessment experiment in a well-controlled laboratory environment under the instructions of ITU-R BT.500-13 [24]. The rendered distorted dynamic digital human videos along with the
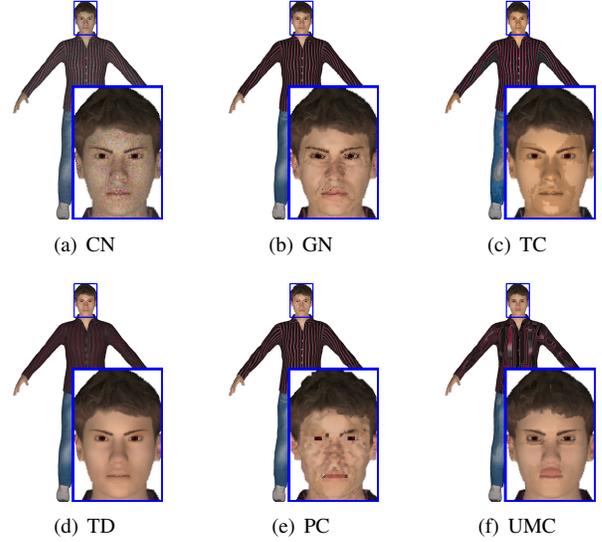
(a) CN     (b) GN     (c) TC

(d) TD     (e) PC     (f) UMC

**Fig. 3**. Projection samples of the model-based distortions, from which we can see that different kinds of distortions can cause diverse perceptual loss to the digital human models.
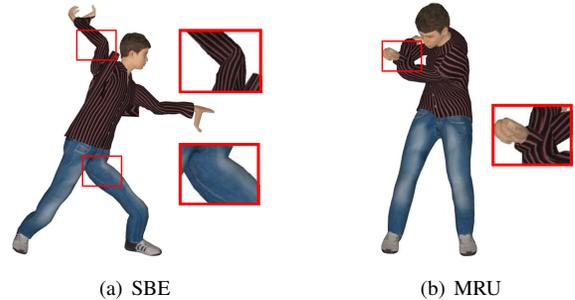


(a) SBE     (b) MRU

**Fig. 4**. Examples of the motion-based distortions. The SBE distortions can severely discord the digital human body and rigidly twist the body joints. The MRU distortions usually cause model clipping, which makes the motion awkward and unnatural.

corresponding reference ones are randomly displayed on a customized graphical subjective quality assessment interface, whose screenshot is shown in Fig. 6. We employ an iMac monitor for display, which supports a resolution up to 4096 × 2304.

A total of 41 subjects (20 males and 21 females) are invited to participate in the subjective experiment. All the subjects are seated from a distance of twice the screen height to the screen in an indoor environment with normal illumination levels. Before the subjective experiment, an instruction session is performed to help the subjects get familiar with the quality assessment task. The whole subjective experiment is split into 16 sessions and each session contains 50 distorted DDH video sequences. There is a 30-minutes
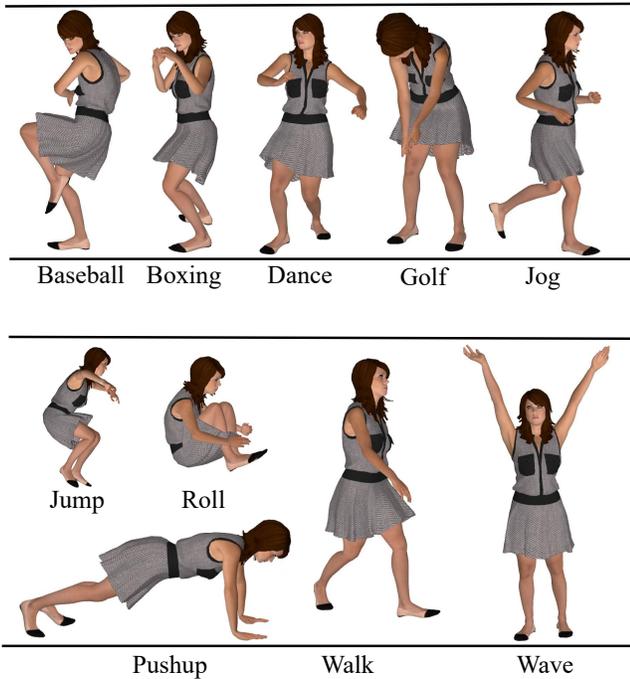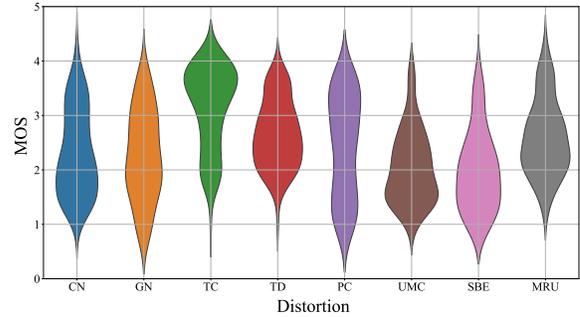
Fig. 5. Illustration of the selected kinds of motion. Simple daily activities such as walk and wave are included. Complicated sports such as baseball and gold are considered as well.
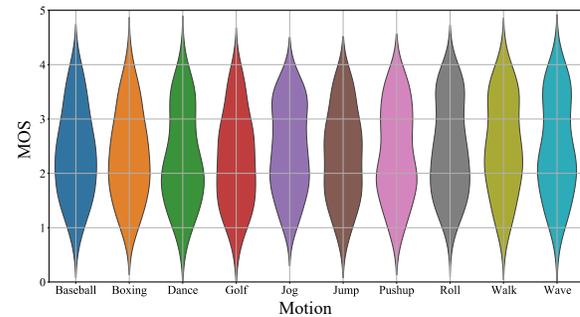


Fig. 6. The screenshot of the subjective quality assessment interface. The reference videos (left) and the distorted videos (right) are displayed at the same time.

break between each session and each subject is allowed to attend no more than 4 sessions in a single day. During each session, the distorted video sequence is played only once and the participants can rate the DDH quality according to the rendered DDH video from 1 to 5, with a minimum interval of 0.1. We ensure that each distorted DDH video is evaluated by the 41 invited participants and 32,800=800×41 subjective ratings are collected in all.



(a)



(b)

Fig. 7. Distributions of the MOSs corresponding to the distortion and motion types.

### III-E. Subjective Data Processing

Following the recommendation of ITU-R BT.500-13 [24], we calculate the z-scores as the quality labels of corresponding DDHs:

$$z_{ij} = \frac{r_{ij} - \mu_i}{\sigma_i}, \tag{1}$$

where $r_{ij}$ represents the quality rating given by the $i$-th subject on the $j$-th DDH, $\mu_i = \frac{1}{N_i} \sum_{j=1}^{N_i} r_{ij}$, $\sigma_i = \sqrt{\frac{1}{N_i - 1} \sum_{j=1}^{N_i} (r_{ij} - \mu_i)}$, and $N_i$ is the number of DDH assessed by subject $i$. Additionally, we reject the quality ratings from unreliable subjects with the recommended subject rejection procedure proposed in [24]. In the end, the z-scores are linearly rescaled to $[1, 5]$ and the mean opinion score (MOS) of DDH $j$ is obtained by averaging the rescaled z-scores:

$$MOS_j = \frac{1}{M} \sum_{i=1}^{M} z'_{ij}, \tag{2}$$

where $MOS_j$ represents the MOS for the $j$-th DDH, $M$ is the number of the valid subjects, and $z'_{ij}$ are the rescaled z-scores.

### III-F. Subjective Data Analysis

We further plot the distributions of the MOSs from the perspective of distortion and motion types, which are shown

in Fig. 7. From the distortion MOS distributions, we can see that the TC distortion tends to have a less negative impact while the SBE distortions seem to cause more severe damage to the visual quality of DDHs. With closer inspections, we can observe that all types of motion exhibit similar MOS distributions, which indicates that the added distortions result in similar perceptual loss regardless of the motion types. Therefore the proposed DDH-QA database can provide useful guidelines for other types of DDH motion.

## IV. BENCHMARK EXPERIMENT

### IV-A. Benchmark Competitors

Since the DDH is usually perceived in the format of animated videos, several state-of-the-art video quality assessment (VQA) methods are employed for validation on the DDH-QA database. The FR methods include PSNR, and SSIM [15], which operate on the frames of DDH videos. The NR methods include BRISQUE [16], NIQE [25], VIIDEO [26], V-BLIINDS [17], TLVQM [18], VIDEVAL [19], VSFA [27], RAPIQUE [20], SimpleVQA [21], and FAST-VQA [22]. Additionally, BRISQUE, NIQE, VIIDEO, V-BLIINDS, TLVQM, and VIDEVAL are handcrafted-based methods while VSFA, RAPIQUE, SimpleVQA, and FAST-VQA are DNN-based methods. It's worth mentioning that we use the source codes provided by the authors and maintain the default setting parameters.

### IV-B. Experimental Setup

The 5-fold cross validation strategy is utilized to train and test the models. Specifically, we split the 10 groups of motion into 5 folds and each fold contains 2 groups of motion. 4 folds are used as the training sets while the left 1 fold is used as the testing set. Such procedure is repeated 5 times so that every fold has been employed as the testing set. The average performance is recorded as the final experimental results. Additionally, for methods that require no training, we simply operate them on the same testing sets and report the average performance.

Four mainstream consistency evaluation criteria are utilized to compare the correlation between the predicted scores and MOSs, which include Spearman Rank Correlation Coefficient (SRCC), Kendall's Rank Correlation Coefficient (KRCC), Pearson Linear Correlation Coefficient (PLCC), and Root Mean Squared Error (RMSE). An excellent model should obtain values of SRCC, KRCC, and PLCC close to 1, and the value of RMSE near 0.

### IV-C. Performance Discussion

The experimental results are shown in Table III, from which we can make several useful conclusions. (a) PSNR and SSIM are the most widely used FR quality assessment metrics in compression and transmission systems. Although they achieve relatively better performance than the NR handcrafted-based methods, they are not effective to deal

**Table III**. Benchmark Performance on the DDH-QA database. Best in bold.

| Ref. | Model | SRCC | PLCC | KRCC | RMSE |
|------|-------|------|------|------|------|
| FR | PSNR | 0.4308 | 0.5458 | 0.3114 | 0.9013 |
| | SSIM | 0.5408 | 0.6057 | 0.3920 | 0.8559 |
| NR | BRISQUE | 0.3664 | 0.4011 | 0.2568 | 1.0067 |
| | NIQE | 0.0923 | 0.2489 | 0.0748 | 1.0418 |
| | VIIDEO | 0.1219 | 0.1829 | 0.0732 | 1.0740 |
| | V-BLIINDS | 0.4807 | 0.4936 | 0.3424 | 0.9564 |
| | TLVQM | 0.2515 | 0.2824 | 0.1729 | 1.0480 |
| | VIDEVAL | 0.2218 | 0.3470 | 0.1622 | 1.0246 |
| | VSFA | 0.5406 | 0.5708 | 0.3858 | 0.9657 |
| | RAPIQUE | 0.1815 | 0.2368 | 0.1246 | 1.0614 |
| | SimpleVQA | **0.7444** | **0.7498** | **0.5452** | **0.7228** |
| | FAST-VQA | 0.5262 | 0.5382 | 0.3657 | 1.0499 |

with practical DDH-QA tasks, which calls for better FR DDH-QA solutions. (b) For the NR-VQA methods, the DNN-based methods except RAPIQUE tend to yield better performance than the handcrafted-based methods. This is because the handcrafted-based methods are based on the natural scene statistics (NSS) prior knowledge, which is learned from natural videos. However, the rendered DDH videos are quite different from the natural videos in both content and distortions, which leads to the ineffectiveness of the handcrafted-based methods. (c) SimpleVQA achieves the highest performance among all the benchmark competitors and is significantly superior to the second-ranking method. We try to give some reasons. Besides using 2D-CNN for spatial feature extraction, SimpleVQA also utilizes 3D-CNN for motion feature extraction, which might be more capable of describing the quality representation of DDHs. To sum up, the existing quality assessment methods still have a long way to go before accurately predicting the visual quality levels of DDHs.

## V. CONCLUSION

In this paper, we propose a large-scale dynamic digital human quality assessment database. One male and one female digital human models are selected as the reference. Then we degrade the reference models with both model-based and motion-based distortions. A total of 800 DDHs are generated and we render the DDHs into 2D animation videos for evaluation. Afterward, we carry out a subjective study to collect the subjective quality judgment for the distorted DDHs. Several state-of-the-art VQA methods are chosen for validation on the proposed DDH-QA database. A comprehensive performance discussion is made as well. We hope our work will draw more attention to the quality assessment of DDHs and inspire future research.

## VI. REFERENCES

[1] Joo H Kim, Karim Abdel-Malek, Jingzhou Yang, and R Timothy Marler, "Prediction and analysis of human motion dynamics performing various tasks," *Inderscience IJHFMS*, vol. 1, no. 1, pp. 69–94, 2006.

[2] Zicheng Zhang, Wei Sun, Yucheng Zhu, Xiongkuo Min, Wu Wei, Ying Chen, and Guangtao Zhai, "Treating point cloud as moving camera videos: A no-reference quality assessment metric," *arXiv preprint arXiv:2208.14085*, 2022.

[3] Zicheng Zhang, Wei Sun, Xiongkuo Min, Quan Zhou, Jun He, Qiyuan Wang, and Guangtao Zhai, "Mm-pcqa: Multi-modal learning for no-reference point cloud quality assessment," *arXiv preprint arXiv:2209.00244*, 2022.

[4] Wenmin Zhu, Xiumin Fan, and Yanxin Zhang, "Applications and research trends of digital human models in the manufacturing industry," *Elsevier VRIH*, vol. 1, no. 6, pp. 558–579, 2019.

[5] Qi Yang, Hao Chen, Zhan Ma, Yiling Xu, Rongjun Tang, and Jun Sun, "Predicting the perceptual quality of point cloud: A 3d-to-2d projection-based exploration," *IEEE TMM*, pp. 1–1, 2020.

[6] Qi Liu, Honglei Su, Zhengfang Duanmu, Wentao Liu, and Zhou Wang, "Perceptual quality assessment of colored 3d point clouds," *IEEE TVCG*, 2022.

[7] Yipeng Liu, Qi Yang, Yiling Xu, and Le Yang, "Point cloud quality assessment: Dataset construction and learning-based no-reference metric," *ACM TOMM*, 2022.

[8] Y. Nehmé, F. Dupont, J. P. Farrugia, P. Le Callet, and G. Lavoué, "Visual quality of 3d meshes with diffuse colors in virtual reality: Subjective and objective evaluation," *IEEE TVCG*, vol. 27, no. 3, pp. 2202–2219, 2021.

[9] Yana Nehmé, Florent Dupont, Jean-Philippe Farrugia, Patrick Le Callet, and Guillaume Lavoué, "Textured mesh quality assessment: Large-scale dataset and deep learning-based quality metric," *arXiv preprint arXiv:2202.02397*, 2022.

[10] Emin Zerman, Cagri Ozcinar, Pan Gao, and Aljosa Smolic, "Textured mesh vs coloured point cloud: A subjective study for volumetric video compression," in *IEEE QoMEX*. IEEE, 2020, pp. 1–6.

[11] Zicheng Zhang, Yingjie Zhou, Wei Sun, Xiongkuo Min, and Guangtao Zhai, "Perceptual quality assessment for digital human heads," *IEEE ICASSP*, 2022.

[12] Zicheng Zhang, Wei Sun, Xiongkuo Min, Tao Wang, Wei Lu, Wenhan Zhu, and Guangtao Zhai, "A no-reference visual quality metric for 3d color meshes," in *IEEE ICMEW*, 2021, pp. 1–6.

[13] Yu Fan, Zicheng Zhang, Wei Sun, Xiongkuo Min, Wei Lu, Tao Wang, Ning Liu, and Guangtao Zhai, "A no-reference quality assessment metric for point cloud based on captured video sequences," *IEEE MMSP*, 2022.

[14] Zicheng Zhang, Wei Sun, Xiongkuo Min, Tao Wang, Wei Lu, and Guangtao Zhai, "No-reference quality assessment for 3d colored point cloud and mesh models," *IEEE TCSVT*, 2022.

[15] Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE TIP*, vol. 13, no. 4, pp. 600–612, 2004.

[16] Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik, "No-reference image quality assessment in the spatial domain," *IEEE TIP*, vol. 21, no. 12, pp. 4695–4708, 2012.

[17] Michele A Saad, Alan C Bovik, and Christophe Charrier, "Blind prediction of natural video quality," *IEEE TIP*, vol. 23, no. 3, pp. 1352–1365, 2014.

[18] Jari Korhonen, "Two-level approach for no-reference consumer video quality assessment," *IEEE TIP*, vol. 28, no. 12, pp. 5923–5938, 2019.

[19] Zhengzhong Tu, Yilin Wang, Neil Birkbeck, Balu Adsumilli, and Alan C Bovik, "Ugc-vqa: Benchmarking blind video quality assessment for user generated content," *IEEE TIP*, vol. 30, pp. 4449–4464, 2021.

[20] Zhengzhong Tu, Xiangxu Yu, Yilin Wang, Neil Birkbeck, Balu Adsumilli, and Alan C Bovik, "Rapique: Rapid and accurate video quality prediction of user generated content," *IEEE DOAJ*, vol. 2, pp. 425–440, 2021.

[21] Wei Sun, Xiongkuo Min, Wei Lu, and Guangtao Zhai, "A deep learning based no-reference quality assessment model for ugc videos," in *ACM MM*, 2022.

[22] Haoning Wu, Chaofeng Chen, Jingwen Hou, Liang Liao, Annan Wang, Wenxiu Sun, Qiong Yan, and Weisi Lin, "Fast-vqa: Efficient end-to-end video quality assessment with fragment sampling," in *ECCV*. 2022, pp. 538–554, Springer.

[23] "Draco 3d data compression," Google, Inc, Oct 1, 2022 [Online],https://github.com/google/draco.

[24] RECOMMENDATION ITU-R BT, "Methodology for the subjective assessment of the quality of television pictures," *International Telecommunication Union*, 2002.

[25] A. Mittal, R. Soundararajan, and A. C. Bovik, "Making a "completely blind" image quality analyzer," *IEEE SPL*, vol. 20, no. 3, pp. 209–212, 2013.

[26] Anish Mittal, Michele A Saad, and Alan C Bovik, "A completely blind video integrity oracle," *IEEE TIP*, vol. 25, no. 1, pp. 289–300, 2015.

[27] Dingquan Li, Tingting Jiang, and Ming Jiang, "Quality assessment of in-the-wild videos," in *ACM MM*, 2019, pp. 2351–2359.