

# LMM-PCQA: Assisting Point Cloud Quality Assessment with LMM

Zicheng Zhang  
zzc1998@sjtu.edu.cn  
Shanghai Jiao Tong University  
Shanghai, China

Wei Sun  
sunguwei@sjtu.edu.cn  
Shanghai Jiao Tong University  
China

Xiaohong Liu  
xiaohongliu@sjtu.edu.cn  
Shanghai Jiao Tong University  
China

Haoning Wu  
haoning001@e.ntu.edu.sg  
Nanyang Technological University  
Singapore

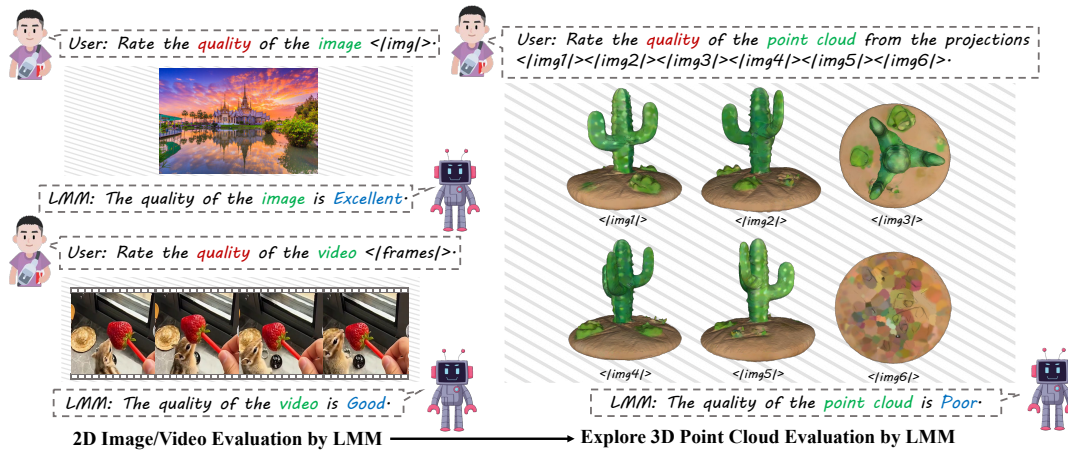
Chaofeng Chen  
chaofeng.chen@e.ntu.edu.sg  
Nanyang Technological University  
Singapore

Weisi Lin  
wslin@e.ntu.edu.sg  
Nanyang Technological University  
Singapore

Yingjie Zhou  
Chunyi Li  
Shanghai Jiao Tong University  
China

Xiongkuo Min<sup>†\*</sup>  
minxiongkuo@sjtu.edu.cn  
Shanghai Jiao Tong University  
China

Guangtao Zhai<sup>†\*</sup>  
zhaiguangtao@sjtu.edu.cn  
Shanghai Jiao Tong University  
China



**Figure 1: Inspired by the impressive quality evaluation ability of LMM on 2D media, we are the first to explore the quality representation potential of LMM on 3D point clouds.**

## ABSTRACT

Although large multi-modality models (LMMs) have seen extensive exploration and application in various quality assessment studies, their integration into Point Cloud Quality Assessment (PCQA) remains unexplored. Given LMMs’ exceptional performance and robustness in low-level vision and quality assessment tasks, this study aims to investigate the feasibility of imparting PCQA knowledge to LMMs through text supervision. To achieve this, we transform quality labels into textual descriptions during the fine-tuning phase,

<sup>†</sup>Corresponding authors<sup>†</sup>.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or to publish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

MM '24, October 28–November 1, 2024, Melbourne, VIC, Australia

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0686-8/24/10

<https://doi.org/10.1145/3664647.3680946>

enabling LMMs to derive quality rating logits from 2D projections of point clouds. To compensate for the loss of perception in the 3D domain, structural features are extracted as well. These quality logits and structural features are then combined and regressed into quality scores. Our experimental results affirm the effectiveness of our approach, showcasing a novel integration of LMMs into PCQA that enhances model understanding and assessment accuracy. We hope our contributions can inspire subsequent investigations into the fusion of LMMs with PCQA, fostering advancements in 3D visual quality analysis and beyond. The code is available at <https://github.com/zzc-1998/LMM-PCQA>.

## CCS CONCEPTS

• **Human-centered computing** → *Visualization design and evaluation methods*; • **Computing methodologies** → **Artificial intelligence**.

## KEYWORDS

Large multi-modality model, Point cloud quality assessment

**ACM Reference Format:**

Zicheng Zhang, Haoning Wu, Yingjie Zhou, Chunyi Li, Wei Sun, Chaofeng Chen, Xiongkuo Min†, Xiaohong Liu, Weisi Lin, and Guangtao Zhai†. 2024. LMM-PCQA: Assisting Point Cloud Quality Assessment with LMM. In *Proceedings of the 32nd ACM International Conference on Multimedia (MM '24)*, October 28–November 1, 2024, Melbourne, VIC, Australia. *Proceedings of the 32nd ACM International Conference on Multimedia (MM'24)*, October 28–November 1, 2024, Melbourne, Australia. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3664647.3680946>

## 1 INTRODUCTION

Point clouds are increasingly used across diverse real-world scenarios, including virtual/augmented reality [15, 21, 31], autonomous vehicles [10], and video post-production [26]. This surge is attributed to their adeptness in three-dimensional representation. Consequently, significant research efforts have been channeled towards enhancing high-level areas like point cloud classification [6, 13, 17, 37, 39, 44, 51], detection [10], and segmentation [8, 24]. Meanwhile, as a key component for ensuring point cloud quality, point cloud quality assessment (PCQA) has seen comparable advancements as well during the last decade. PCQA's objective is to appraise point clouds' visual quality, a pivotal factor for refining simplification and compression strategies in practical applications [11], and to elevate the Quality of Experience (QoE) for end-users. Generally, PCQA methods are divided into three categories, depending on their reliance on reference point clouds: Full-Reference PCQA (FR-PCQA), Reduced-Reference PCQA (RR-PCQA), and No-Reference PCQA (NR-PCQA) respective. The availability of pristine reference point clouds is often limited in real-world applications, emphasizing the demand for NR-PCQA approaches, which is why our research primarily focuses on NR-PCQA.

There are already some cutting-edge studies that have begun applying Large Multi-Modality Models (LMMs) to low-level vision and quality assessment fields [16, 40–42, 52], achieving notable success. LMMs demonstrate highly competitive performance and robustness in these tasks. However, the main focus of these studies has remained on two-dimensional (2D) media such as images and videos, **while no research has explored the possibility of applying LMMs to three-dimensional (3D) media like point clouds**. It is known that both 2D and 3D media exhibit similar distortions, i.e., blur and noise. Given the robust quality perception of LMMs in 2D, **we can hold the hypothesis that LMM also has significant quality perception abilities in 3D point clouds**. Hence, investigating the application of LMMs for point cloud quality assessment is not only valuable but also meaningful, reflecting their established visual perception strengths. Therefore, in this study, we carry out a novel method named **LMM-PCQA** to provide an interesting solution for handling the PCQA problem with the assistance of LMM.

First, **we treat the point clouds as sequences of projections** to enable LMM to perceive the point cloud visual quality. Afterward, we try to **teach LMM about the quality alignment** between the predefined 5-level **qualitative adjectives** (i.e., *excellent, good, fair, poor, bad*) and point cloud projections. Specifically, we employ the existing PCQA databases to provide the necessary knowledge, where the quality labels are transformed into corresponding

**qualitative adjectives**. Then we specially design a prompt structure to produce the question-answer pairs, which are composed of the question *'Rate the quality of the point cloud from the projections [img1],[img2],[img3],[img4],[img5],[img6]'* and the answer *'The quality of the point cloud is excellent/good/fair/poor/bad'*. These question-answer pairs can be utilized to teach LMM PCQA knowledge during the fine-tuning stage. After the instruction tuning, we can expect the trained LMM to give the predicted `[SCORE_TOKEN]` with the same prompt structure (the **qualitative adjectives** position is left blank for LMM response), which is shown in Fig. 2. The predicted `[SCORE_TOKEN]` can be recognized as a probability map to the **qualitative adjectives**, and we convert it into 5-level probabilities as the LMM evaluation results.

Secondly, to address the potential insensitivity to geometric distortions (e.g., compression, downsampling) when only analyzing projections, **we propose the extraction of multi-scale structural features**. This approach enhances the LMM-PCQA's holistic comprehension of point cloud visual quality. The point clouds are converted into quality-aware structural domains, a technique validated for effective quality feature extraction in prior research [4, 58, 65]. We modify the scale parameters in the k-nearest neighbors (k-NN) algorithm to offer a multi-scale perspective, aligning with the human vision system's perception mechanism. Subsequently, key statistical parameters are used to quantify structural distortions within these domains. Finally, we combine the LMM's evaluative results and the structural features, utilizing support vector regression (SVR) to derive the quality values. The experimental outcomes affirm that our LMM-PCQA model is on par with, or surpasses, current leading PCQA methods.

By carrying out LMM-PCQA, the contributions of this paper can be summarized as follows:

- **We are the first to employ LMM for PCQA tasks.** We design a novel prompt structure to enable the LMM to perceive the point cloud visual quality. The existing PCQA databases are converted into question-answer pairs, which are then used to **inject PCQA knowledge into LMM**.
- **We propose the extraction of multi-scale structural features.** By processing point clouds into multi-scale domains, we quantify the geometry distortions via key statistic parameters estimation, which helps LMM gain a more comprehensive understanding of point cloud visual quality.
- **LMM-PCQA demonstrates exceptional performance across various PCQA databases.** The ablation study and cross-database evaluations further validate the logical design of LMM-PCQA and its robust generalization capabilities.

## 2 RELATED WORKS

### 2.1 LMM for Quality Assessment

Recent studies are exploring LMMs in visual quality assessment. X-IQE [7] uses them to evaluate text-to-image generation via a chain of thoughts strategy. Q-Bench [40] introduces a binary softmax method for LMMs to produce quantifiable scores by assessing tokens *good* and *poor*. Q-Instruct [41] improves this by fine-tuning LMMs through text-based questioning for specific visual queries. Q-Align [42] adopts human-like evaluation mechanisms, enhancing visual quality scoring.

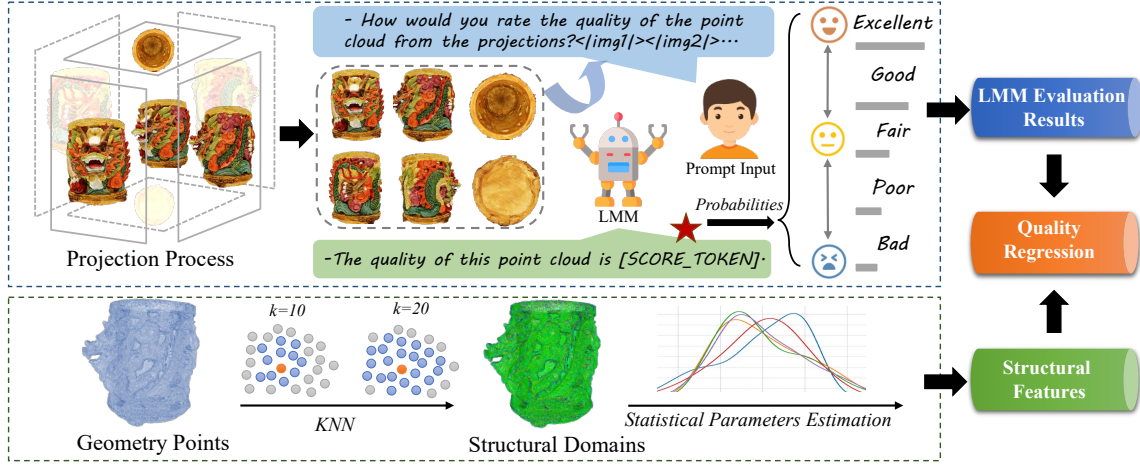


Figure 2: The framework of the proposed method.

## 2.2 PCQA Development

In the early stages of PCQA, the MPEG group developed Full-Reference PCQA (FR-PCQA) methods like p2point [27] and p2plane [34] for point cloud quality assessment, and later introduced a color-based PSNR-yuv method [35]. Despite challenges in addressing complex distortions, advanced FR-PCQA metrics like PCQM [28], GraphSIM [48], and PointSSIM [3] were developed to integrate structural features for enhanced performance.

Building on progress in no-reference image and video quality assessment (NR-I/VQA) [14, 23, 53, 54, 56, 61], various NR-PCQA methods have evolved. Techniques include using CNNs for quality regression with handcrafted features [9], multi-view projections for feature extraction [22, 58], and direct prediction with sparse CNNs [25]. Innovative approaches also involve converting point clouds to videos for VQA [12, 59], employing structure-guided resampling [64], integrating multi-projections for simpler feature extraction [62], adapting image quality assessment methods for point cloud rendering [47], and using non-local geometry and color gradient models for quality estimation [38]. Recent frameworks like MM-PCQA [60] and pmBQA [45] leverage multi-modal learning for improved PCQA. Emerging research explores text modality’s potential in quality assessment learning.

## 3 PROPOSED METHOD

The framework of the proposed method is briefly illustrated in Fig. 2, which includes the LMM evaluation module, the structural feature extraction module, and the quality regression module.

### 3.1 LMM Evaluation

**3.1.1 Projection Acquisition.** Consider a colored point cloud  $\mathbf{P} = (p_i, c_i)_{i=1}^N$ , where  $p_i \in \mathbb{R}^{1 \times 3}$  represents the single point consisting of geometry coordinates,  $c_i \in \mathbb{R}^{1 \times 3}$  represents the RGB color attributes, and  $N$  denotes the total count of points. We then adopt the conventional cube-like viewpoints configuration, which is inspired by the projections approach discussed in [2, 35]. As illustrated in Fig. 2, six orthogonal viewpoints are utilized, each mapping to one of the cube’s six faces for generating the projections. For a point cloud  $\mathbf{P}$ , the rendering process is derived as:

$$\begin{aligned} \mathbf{I} &= \psi(\mathbf{P}), \\ \mathbf{I} &= \{\mathcal{I}_k | k = 1, \dots, 6\}, \end{aligned} \quad (1)$$

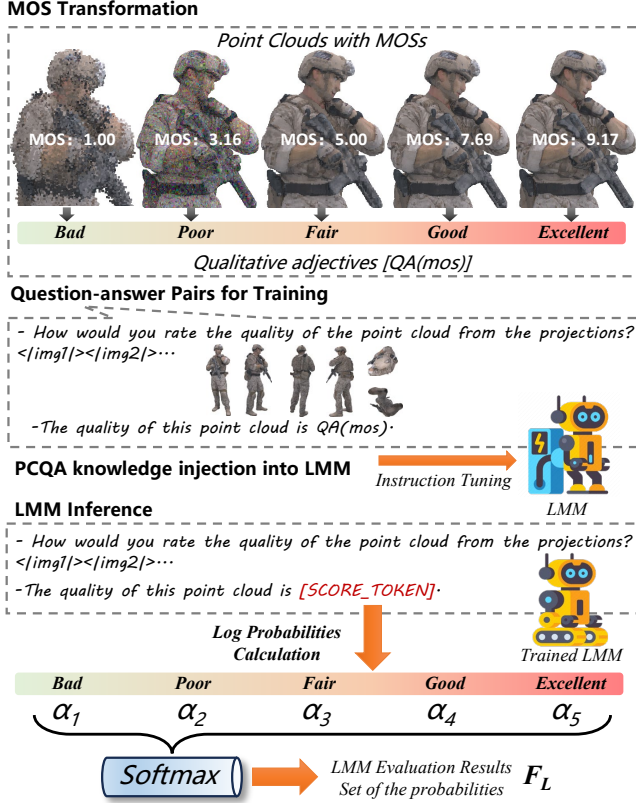
where  $\mathbf{I}$  represents the set of the 6 rendered projections and  $\psi(\cdot)$  denotes the Open3D-based [63] projections capturing process.

**3.1.2 How to inject PCQA knowledge into LMMs?** LMM has shown competitive performance in 2D image/video quality assessment [40–42], therefore it is feasible to apply LMM for PCQA tasks by taking the point cloud as a sequence of projections. Then it is vital to solve the core problem of **PCQA Knowledge Injection**. Following the common approaches of training LMMs [18, 49], it is natural to come up with the solution of instructing the LMM with question-answer pairs regarding PCQA issues. Thus we carry out the specific prompt structure as follows:

-How would you rate the quality of the point cloud from the projections?<img1><img2>...  
-The quality of this point cloud is [QA(mos)].

where <img1><img2>... stands for the image set of projections and QA(s) is the **qualitative adjective** of the point cloud which can be obtained from the mean opinion score (MOS) corresponding to the point cloud. Afterward, we can use the designed question-answer pairs to teach LMM PCQA knowledge.

**3.1.3 Transformation from MOS to quality rating.** In daily experiences, humans often give feedback using **qualitative adjectives** (such as good, bad, superb) instead of **numerical ratings** (like 9.2, 2.5, 7.1). Hence, implementing visual scoring activities with level-based ratings taps into this *natural tendency* of humans (to offer qualitative adjectives). Similarly, the perception and expression of LMM are akin to humans, which have a better understanding and perception of **qualitative adjectives**. Therefore, converting MOS into corresponding **qualitative adjectives** for its learning is more intuitive than having it learn directly from numbers. Specifically, the transformation of MOS can be achieved by evenly splitting the range from the highest score (M) to the lowest score (m) into



**Figure 3: Illustration of the LLM evaluation pipeline. The point clouds with MOSs are transformed into question-answer pairs for LLM tuning. The LLM evaluation results can be obtained as the set of the probabilities to the predefined qualitative adjectives.**

five unique intervals, with scores in each interval designated as corresponding quality levels:

$$QA(mos) = w_i \text{ if } m + \frac{i-1}{5} \times (M-m) < mos \leq m + \frac{i}{5} \times (M-m), \quad (2)$$

where  $\{w_i\}_{i=1}^5 = \{bad, poor, fair, good, excellent\}$  are the standard text rating levels as defined by ITU [1].

**3.1.4 Obtaining evaluation results via LLM inference.** After training, we can get the evaluation results with the same prompt structure and get the response  $[SCORE\_TOKEN]$  via LLM inference. The  $[SCORE\_TOKEN]$  can be recognized as a log probability map to the **qualitative adjectives**. Then we can compute the final probabilities to the 5-level **qualitative adjectives** from the corresponding log probabilities via softmax as the LLM evaluation results:

$$F_L = \left\{ \frac{e^{\alpha_i}}{\sum_{j=1}^5 e^{\alpha_j}} \right\}_{i=1}^5, \quad (3)$$

where  $\alpha_i$  indicates the log probability of  $i$ -th **qualitative adjectives** and  $F_L$  represents the LLM evaluation results which consist of 5 probabilities after softmax.

## 3.2 Structural Feature Extraction

The inadequacy of single-modality information from projections for point cloud quality evaluation has been established [45, 60]. Therefore, our approach enhances the accuracy of LMM evaluation results (projections only) by integrating geometric structural features, aiming for a more detailed and accurate assessment.

**3.2.1 Structural Domain.** Given the point cloud  $P = (p_i, c_i)_{i=1}^N$ , the neighborhood  $P_{Nbi}$  of each point  $p_i$  can be obtained utilizing the  $k$ -nearest neighbors ( $k$ -NN) algorithm:

$$P_{Nb} = \text{KNN}(P),$$

$$\text{Dist}(p, q) = \sqrt{(p_x - q_x)^2 + (p_y - q_y)^2 + (p_z - q_z)^2}, \quad (4)$$

where  $N$  represents the total count of points within the point cloud, the term  $P_{Nb}$  denotes the collection of neighborhoods, while  $\text{KNN}(\cdot)$  signifies the function of the  $k$ -nearest neighbors algorithm. The distance between points  $p$  and  $q$  is calculated using the Euclidean distance, expressed as  $\text{Dist}(p, q)$ . Given the neighborhood set  $P_{Nbi}$  of point  $p_i$ , we can define the covariance matrix  $C_i$  for each point  $p_i$ , which is characterized by its 3D geometric coordinates:

$$C_i = \frac{1}{K} \sum_{j=1}^K (p_{n_j} - \hat{p})(p_{n_j} - \hat{p})^\top, \quad (5)$$

$$\{p_{n_1}, \dots, p_{n_K}\} \in P_{Nbi},$$

where the term  $K$  denotes size of neighborhood set  $P_{Nbi}$ ,  $p_{n_j}$  is the  $j$ -th neighboring point in  $P_{Nbi}$ ,  $\hat{p}$  is the centroid of this neighborhood,  $p_{n_j}$  and  $\hat{p}$  are vectors with dimensions  $\mathbb{R}^{3 \times 1}$ , while  $C_i$  is a matrix with dimensions  $\mathbb{R}^{3 \times 3}$ . Consequently, the eigenvectors for the covariance matrix  $C_i$  can be derived as follows:

$$C_i \cdot v_l = \lambda_l \cdot v_l, l \in \{1, 2, 3\}, \quad (6)$$

where  $(\lambda_1, \lambda_2, \lambda_3)$  stand for the eigenvalues and  $(v_1, v_2, v_3)$  represent the respective eigenvectors, with  $\lambda_1 > \lambda_2 > \lambda_3$ . Consequently, we derive three eigenvalues for each point  $p_i$  within the point cloud  $P$ . Then we can compute the *linearity* and *planarity* structural domains of the point cloud as:

$$\text{Lin}(p_i) = \frac{\lambda_1 - \lambda_2}{\lambda_1}, \quad (7)$$

$$\text{Pla}(p_i) = \frac{\lambda_2 - \lambda_3}{\lambda_1},$$

where  $\text{Lin}(p_i)$  and  $\text{Pla}(p_i)$  represent the *linearity* and *planarity* values for point  $p_i$ . The chosen structural domains (*linearity*, *planarity*) have been demonstrated to exhibit a strong correlation with geometric visual losses, such as compression and downsampling, and have been extensively utilized in numerous PCQA tasks [4, 58, 65].

**3.2.2 Multi-scale Perception.** The multi-scale nature of point cloud visual perception has been noted in the literature [57]. To account for this, we compute structural domains across various scales by varying the scale parameter  $k$  in the  $\text{KNN}(\cdot)$  process. This approach allows us to derive the multi-scale structural domains as follows:

$$\mathcal{D}_{k=k_{ms}} = \mathcal{S}(\text{KNN}(P)), \mathcal{D} \in \{\text{Lin}, \text{Pla}\}, \quad (8)$$

where  $\mathcal{D}_{k=k_{ms}}$  denotes the multi-scale structural domains, with  $k_{ms}$  represents the set of scale parameters, and  $\mathcal{S}(\cdot)$  refers to the process of calculating structural domains as previously described.



In our study, we establish the default set of scale parameters as  $\{10, 20\}$ , signifying that the *linearity* and *planarity* domains are computed using the 10-nearest-neighbor and 20-nearest-neighbor configurations, respectively.

**3.2.3 Statistical Parameters Estimation.** To quantify the quality representation from the structural domains, we employ some basic statistical parameters estimation process:

$$F_S = \{avg(\mathcal{D}_{k=k_{ms}}), std(\mathcal{D}_{k=k_{ms}}), ent(\mathcal{D}_{k=k_{ms}})\}, \quad (9)$$

$$\mathcal{D} \in \{Lin, Pla\},$$

where  $avg(\cdot)$ ,  $std(\cdot)$ , and  $ent(\cdot)$  represent the average function, standard deviation function, and entropy function respectively, and  $F_S$  indicates the set of the final extracted structural features.

### 3.3 Quality Regression

To clearly demonstrate the efficacy of the proposed features, we merge the LMM evaluation results with the structural features, and then incorporate them into the visual quality score using support vector regression (SVR):

$$Q = SVR(F_L \oplus F_S), \quad (10)$$

where  $Q$  indicates the quality values,  $SVR(\cdot)$  represents the SVR regression process, and  $\oplus$  denotes the concatenation process.

## 4 EXPERIMENT

### 4.1 Validation Databases

To assess the efficacy of the proposed method, we employ the SJTU-PCQA database [46], the Waterloo point cloud assessment database (WPC)[19], and the WPC2.0 database [20] for validation. The SJTU-PCQA database contains 9 reference point clouds, subjected to seven distortion types (compression, color noise, geometric shift, down-sampling, and three mixed distortions) at six levels, yielding a total of 378 ( $9 \times 7 \times 6$ ) distorted point clouds. The WPC database includes 20 reference point clouds, each modified by four distortions (down-sampling, Gaussian white noise, Geometry-based Point Cloud Compression (G-PCC), and Video-based Point Cloud Compression (V-PCC)), resulting in 740 ( $20 \times 37$ ) distorted point clouds. Meanwhile, the WPC2.0 database features 16 reference point clouds, each undergoing 25 different V-PCC degradation settings, leading to 400 ( $16 \times 25$ ) distorted point clouds.

### 4.2 Implementation Details

**4.2.1 LMM Training.** Following the mainstream choice of LMM-involved quality assessment methods [41–43], we select the mPLUG-Owl-2 [50] as the LMM model in this paper. The model comprises a CLIP-ViT-Large [32] visual encoder  $\mathcal{E}_v$  with 304 million parameters, a visual abstractor  $\hat{\mathcal{E}}_v$  with 82 million parameters, and the LLaMA2-7B [36] LLM  $\mathcal{L}$  on top of the visual modules, which integrates an additional multi-way module from mPLUG-Owl2, totaling 7.8 billion parameters. Input projections are initially squared through padding before being resized to  $448 \times 448$ . Let  $\mathcal{E}_t$  represent the text embedding layer, with input projections denoted as  $\langle img1 \rangle \langle img2 \rangle \dots$  and the text prompt as  $t$ , the detailed formulation of the used LMM

model can be expressed as follows:

$$\begin{aligned} \mathcal{H}_v &= \hat{\mathcal{E}}_v(\mathcal{E}_v(\langle img1 \rangle \langle img2 \rangle \dots)), \\ \mathcal{H}_t &= \mathcal{E}_t(t), \\ \mathcal{H} &= \mathcal{H}_v \oplus \mathcal{H}_t, \\ \mathcal{O} &= \mathcal{L}(\mathcal{H}), \end{aligned} \quad (11)$$

where  $\mathcal{H}_v$  and  $\mathcal{H}_t$  represent the abstracted tokens for the visual and text input respectively,  $\mathcal{O}$  stands for the output. For all PCQA databases, the batch size is maintained at 64. The learning rate is fixed at  $2 \times 10^{-5}$ , with the training process extending over 2 epochs for each variant. We utilize the common GPT [33] loss mechanism, specifically the cross-entropy between the predicted logits and actual labels. The evaluation of performance metrics is conducted using the final weights obtained post-training. Four NVIDIA A100 80G GPUs are employed for the training phase, while a single RTX3090 24G GPU is used to measure inference latency. In the inference stage, only the input texts preceding the  $[SCORE\_TOKEN]$  are inputted into the LMM, leading to the final element of  $\mathcal{O}$  representing the targeted probability map.

**4.2.2 Validation Strategy.** Following the methodologies in [5, 12, 60], we utilize a k-fold cross-validation approach in our experiments to ensure a dependable performance evaluation of our proposed method. The SJTU-PCQA, WPC, and WPC2.0 databases consist of 9, 20, and 16 point cloud groups, respectively, leading us to adopt 9-fold, 5-fold, and 4-fold cross-validation for these databases to achieve an approximate 8:2 train-test split. The average of the performance metrics is considered the definitive result. It is crucial to note that the training and testing sets are mutually exclusive to prevent content overlap. For FR-PCQA methods, which do not require training, we evaluate them using the same test sets and report the average performance.

### 4.3 Competitors

17 state-of-the-art quality assessment methods are selected for comparison, which consists of 8 FR-PCQA and 9 NR-PCQA methods:

- The FR-PCQA methods include MSE-p2point (MSE-p2po) [27], Hausdorff-p2point (HD-p2po) [27], MSE-p2plane (MSE-p2pl) [34], Hausdorff-p2plane (HD-p2pl) [34], PSNR-yuv [35], PCQM [28], GraphSIM [48], and PointSSIM [3].
- The NR-PCQA methods include BRISQUE [29], NIQE [30], IL-NIQE [55], IT-PCQA [47], ResSCNN [25], PQA-net [22], 3D-NSS [58], GMS-3DQA [62] and MM-PCQA [60].

Note that BRISQUE, NIQE, IL-NIQE are image-based quality assessment metrics and are validated on the same projections.

### 4.4 Evaluation Criteria

Four mainstream evaluation criteria in the quality assessment field 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), Root Mean Squared Error (RMSE). An excellent quality assessment model should obtain values of SRCC, KRCC, PLCC close to 1 and RMSE to 0.

**Table 1: Performance on the SJTU-PCQA, WPC, and WPC2.0 databases. Best in red, second in blue.**

Type	Methods	SJTU-PCQA				WPC				WPC2.0			
		SRCC↑	PLCC↑	KRCC↑	RMSE↓	SRCC↑	PLCC↑	KRCC↑	RMSE↓	SRCC↑	PLCC↑	KRCC↑	RMSE↓
FR	MSE-p2po	0.7294	0.8123	0.5617	1.3613	0.4558	0.4852	0.3182	19.8943	0.4315	0.4626	0.3082	19.1605
	HD-p2po	0.7157	0.7753	0.5447	1.4475	0.2786	0.3972	0.1943	20.8990	0.3587	0.4561	0.2641	18.8976
	MSE-p2pl	0.6277	0.5940	0.4825	2.2815	0.3281	0.2695	0.2249	22.8226	0.4136	0.4104	0.2965	21.0400
	HD-p2pl	0.6441	0.6874	0.4565	2.1255	0.2827	0.2753	0.1696	21.9893	0.4074	0.4402	0.3174	19.5154
	PSNR-yuv	0.7950	0.8170	0.6196	1.3151	0.4493	0.5304	0.3198	19.3119	0.3732	0.3557	0.2277	20.1465
	PCQM	0.8644	0.8853	0.7086	1.0862	0.7434	0.7499	0.5601	15.1639	0.6825	0.6923	0.4929	15.6314
	GraphSIM	0.8783	0.8449	0.6947	1.0321	0.5831	0.6163	0.4194	17.1939	0.7405	0.7512	0.5533	14.9922
	PointSSIM	0.6867	0.7136	0.4964	1.7001	0.4542	0.4667	0.3278	20.2733	0.4810	0.4705	0.2978	19.3917
NR	BRISQUE	0.3975	0.4214	0.2966	2.0937	0.2614	0.3155	0.2088	21.1736	0.0820	0.3353	0.0487	21.6679
	NIQE	0.1379	0.2420	0.1009	2.2622	0.1136	0.2225	0.0953	23.1415	0.1865	0.2925	0.1335	22.5146
	IL-NIQE	0.0837	0.1603	0.0594	2.3378	0.0913	0.1422	0.0853	24.0133	0.0911	0.1233	0.0714	23.9987
	IT-PCQA	0.8651	0.8283	0.6430	1.1661	0.4870	0.4329	0.3006	19.8960	0.5661	0.5432	0.3477	18.7224
	ResSCNN	0.8600	0.8100	-	-	-	-	-	-	0.7500	0.7200	-	-
	PQA-net	0.8372	0.8586	0.6304	1.0719	0.7026	0.7122	0.4939	15.0812	0.6191	0.6426	0.4606	16.9756
	3D-NSS	0.7144	0.7382	0.5174	1.7686	0.6479	0.6514	0.4417	16.5716	0.5077	0.5699	0.3638	17.7219
	GMS-3DQA	<b>0.9108</b>	0.9177	0.7735	0.7872	0.8308	0.8338	0.6457	<b>12.2292</b>	<b>0.8272</b>	<b>0.8218</b>	<b>0.6277</b>	<b>12.9904</b>
	MM-PCQA	0.9103	<b>0.9226</b>	<b>0.7838</b>	<b>0.7716</b>	<b>0.8414</b>	<b>0.8556</b>	<b>0.6513</b>	12.3506	0.8023	0.8024	0.6202	13.4289
	LMM-PCQA(Ours)	<b>0.9376</b>	<b>0.9404</b>	<b>0.8002</b>	<b>0.7175</b>	<b>0.8825</b>	<b>0.8739</b>	<b>0.7064</b>	<b>11.8171</b>	<b>0.8614</b>	<b>0.8634</b>	<b>0.6723</b>	<b>10.6924</b>

## 4.5 Performance Discussion

The performance comparison between the proposed LMM-PCQA and other PCQA competitors on the SJTU-PCQA, WPC, and WPC2.0 databases is illustrated in Table 1, from which we can draw several conclusions: 1) The LMM-PCQA demonstrates superior performance across all three PCQA databases, even outperforming FR-PCQA methods. For instance, the proposed LMM-PCQA surpasses the second-best NR-PCQA method by about 0.027 (against GMS-3DQA), 0.041 (against MM-PCQA), and 0.034 (against GMS-3DQA) on the SJTU-PCQA, WPC, and WPC2.0 databases from the SRCC values. This highlights LMM’s capability to effectively assimilate PCQA knowledge and apply it actively. The consistency in LMM-PCQA’s performance across various databases underscores its potential to set new baselines in the PCQA field. 2) All PCQA competitors generally perform better on the SJTU-PCQA database but face significant performance declines on the WPC and WPC2.0 databases. In contrast, LMM-PCQA exhibits the smallest performance drops, with decreases of approximately 0.06 and 0.08 in SRCC values when transitioning from the SJTU-PCQA to the WPC and WPC2.0 databases, respectively. This performance stability underscores LMM-PCQA’s robustness and superior ability to handle diverse content effectively, showcasing its adaptability and consistency in quality assessment across different point cloud databases.

## 4.6 Ablation Study

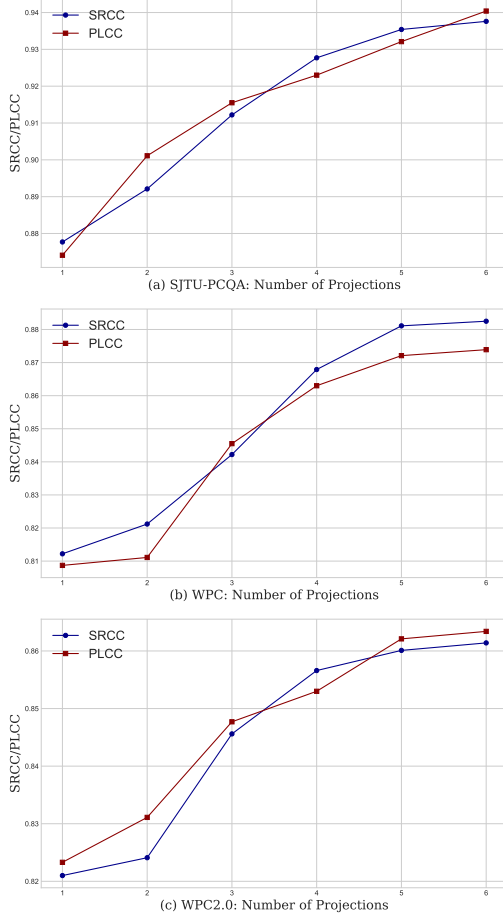
**4.6.1 Contributions of LMM evaluation results and structural features.** To fully investigate the contributions and validate the rationality behind the proposed dual streams of features, we decide to undertake an ablation study in this section. The results, as detailed in Table 2, clearly demonstrate that the integration of both feature streams leads to superior performance compared to employing a single feature stream. Upon a detailed examination, it is apparent

**Table 2: Contributions of LMM evaluation results and structural features, where ‘w/o LMM’ indicates excluding the LMM evaluation results, ‘w/o Structural’ indicates excluding the structural features.**

Modal	SJTU-PCQA		WPC		WPC2.0	
	SRCC↑	PLCC↑	SRCC↑	PLCC↑	SRCC↑	PLCC↑
w/o LMM	0.6650	0.7274	0.3598	0.3523	0.3847	0.3951
w/o Structural	<b>0.9081</b>	<b>0.9158</b>	<b>0.8488</b>	<b>0.8271</b>	<b>0.8258</b>	<b>0.8381</b>
LMM + Structural	<b>0.9376</b>	<b>0.9404</b>	<b>0.8825</b>	<b>0.8739</b>	<b>0.8614</b>	<b>0.8634</b>

that LMM evaluation results markedly surpass the structural features in terms of performance. This disparity primarily stems from the fact that some distorted point clouds within the SJTU-PCQA, WPC, and WPC2.0 databases suffer only from color distortions. Structural features, derived from geometric information, inherently lack the capability to recognize these color distortions, resulting in their relatively lower performance in total. However, incorporating structural features enhances the comprehensive understanding of point cloud quality, contributing to the refinement and elevation of the assessment precision in LMM evaluation results. This combined method highlights how important it is to look at features from different angles to fully understand and capture the subtle details of point cloud quality.

**4.6.2 Influence of the number of projections.** The proposed LMM-PCQA utilizes 6 projections as default. In this section, we further change the number of used projections to test the corresponding performance influence. Specifically, we randomly select 1-6 projections from the cube-like 6 projections setting as the input projections of LMM-PCQA. The performance tendency is illustrated in Fig. 4, from which we can draw several interesting conclusions: 1)



**Figure 4: SRCC/PLCC performance tendency according to the number of used projections on the SJTU-PCQA, WPC, and WPC2.0 databases.**

As the number of projections increases, the performance of LMM-PCQA also improves correspondingly, indicating that increasing the number of projections can encompass more effective quality information, thereby enhancing the final performance. 2) Specifically, when the number of projections increases from 2 to 5, the improvement is significantly more pronounced. This suggests that at this stage, the quality information is not yet redundant, and the benefit of increasing the number of projections is relatively large. However, when the number of projections increases from 5 to 6, the performance improvement is relatively low, indicating that the quality information has become somewhat saturated, and further increases in the number of projections yield diminishing returns.

**4.6.3 Effect of the multi-scale structural features.** To quantify the contributions of the multi-scale mechanism, we validate the performance of structural features with different scale parameters under two settings: with LMM evaluation results and without LMM evaluation results. The experimental performance is listed in Table 3. From the table, we can find that the multi-scale structural features with  $k=10,20$  perform better than the single-scale features whether

**Table 3: Performance of the multi-scale structural features, where  $k$  is the scale parameter of the KNN algorithm.**

Model	SJTU-PCQA		WPC		WPC2.0	
	SRCC $\uparrow$	PLCC $\uparrow$	SRCC $\uparrow$	PLCC $\uparrow$	SRCC $\uparrow$	PLCC $\uparrow$
w/o LMM						
$k=10$	<b>0.6090</b>	<b>0.6584</b>	<b>0.3261</b>	<b>0.3482</b>	<b>0.3366</b>	0.2777
$k=20$	0.5920	0.6311	0.1795	0.2720	0.3224	<b>0.3129</b>
$k=10,20$	<b>0.6650</b>	<b>0.7274</b>	<b>0.3598</b>	<b>0.3523</b>	<b>0.3847</b>	<b>0.3951</b>
with LMM						
$k=10$	0.9140	0.9176	<b>0.8564</b>	<b>0.8554</b>	0.8432	0.8466
$k=20$	<b>0.9199</b>	<b>0.9179</b>	0.8466	0.8488	<b>0.8578</b>	<b>0.8562</b>
$k=10,20$	<b>0.9376</b>	<b>0.9404</b>	<b>0.8825</b>	<b>0.8739</b>	<b>0.8614</b>	<b>0.8634</b>

**Table 4: The cross-database evaluation performance, ‘WPC $\rightarrow$ SJTU-PCQA’ signifies that the model is trained using the WPC database and tested according to the standard testing protocol of the SJTU-PCQA database. We eliminate those point cloud groups from the WPC database that have reference counterparts in the WPC2.0 testing sets, thereby preventing content duplication.**

Model	WPC $\rightarrow$ SJTU-PCQA		WPC $\rightarrow$ WPC2.0	
	SRCC $\uparrow$	PLCC $\uparrow$	SRCC $\uparrow$	PLCC $\uparrow$
PQA-net	0.5411	0.6102	0.6006	0.6377
3D-NSS	0.1817	0.2344	0.4933	0.5613
GMS-3DQA	0.7421	0.7611	0.7822	0.7714
MM-PCQA	<b>0.7991</b>	<b>0.7902</b>	<b>0.7917</b>	<b>0.7935</b>
LMM-PCQA(Ours)	<b>0.8246</b>	<b>0.7999</b>	<b>0.8385</b>	<b>0.8387</b>

the LMM evaluation results are involved or not, which confirms the effectiveness of the proposed multi-scale mechanism. This can be attributed to that humans tend to perceive the visual quality of point clouds from a multi-scale perspective.

## 4.7 Cross-database Validation

To the generalization ability of the proposed LMM-PCQA, we conduct the cross-database validation in this section. Considering that the SJTU-PCQA, WPC, and WPC2.0 databases contain 378, 740, and 400 distorted point clouds respectively, we pre-train LMM-PCQA on the WPC database (largest in scale) and validate the performance on the SJTU-PCQA and WPC2.0 databases (smaller in scale). The competitive NR-PCQA methods (PQA-net, 3d-NSS, GMS-3DQA, and MM-PCQA) are included for comparison. The experimental performance is shown in Table 4, from which we can make several observations: 1) The proposed LMM-PCQA achieves the best cross-database validation performance against all competitors, which confirms the strong generalization ability of LMM-PCQA. 2) Most methods obtain higher WPC $\rightarrow$ WPC2.0 performance than WPC $\rightarrow$ SJTU-PCQA performance. This might be because the WPC2.0 database contains only compression distortions, which undergo part of a similar distortion generation process of the WPC database. Therefore, the quality representation learned from the WPC database is more effective on the WPC2.0 database.

**Table 5: Distortion-specific performance results on the SJTU-PCQA database, where OT represents octree-based compression, CN represents color noise, DS represents down-sampling, DS+CN represents down-sampling and color noise, DS+GGN represents down-sampling and geometry Gaussian noise, GGN represents geometry Gaussian noise, and CN+GGN represents color noise and geometry Gaussian noise respectively.**

Distortion	OT		CN		DS		DS+CN		DS+GGN		GGN		CN+GGN	
	SR↑	PL↑	SR↑	PL↑	SR↑	PL↑	SR↑	PL↑	SR↑	PL↑	SR↑	PL↑	SR↑	PL↑
MSE-p2po	0.71	0.76	nan	nan	0.93	0.95	<b>0.96</b>	0.87	0.96	0.89	<b>0.98</b>	0.89	<b>0.99</b>	0.90
HD-p2po	0.64	0.69	nan	nan	0.82	0.88	0.80	0.75	0.94	0.91	<b>0.98</b>	0.91	<b>0.99</b>	0.91
MSE-p2pl	0.55	0.62	nan	nan	0.87	0.92	0.85	0.81	0.96	0.75	<b>0.97</b>	0.85	<b>0.98</b>	0.86
HD-p2pl	0.54	0.58	nan	nan	0.82	0.87	0.81	0.79	0.94	0.77	0.95	0.88	0.97	0.85
PSNR-yuv	0.59	0.54	0.86	0.87	0.91	0.91	<b>0.96</b>	0.91	<b>0.97</b>	0.94	<b>0.98</b>	0.95	<b>0.99</b>	0.96
PCQM	0.80	0.84	0.86	0.85	0.93	0.96	<b>0.97</b>	<b>0.94</b>	0.96	0.90	<b>0.98</b>	0.93	<b>0.99</b>	0.93
GraphSIM	0.82	0.81	0.82	0.90	<b>0.96</b>	<b>0.97</b>	0.91	<b>0.95</b>	0.95	0.95	0.96	<b>0.97</b>	0.97	<b>0.98</b>
PointSSIM	0.80	<b>0.88</b>	0.87	0.87	0.93	0.93	<b>0.97</b>	0.93	0.96	<b>0.97</b>	<b>0.98</b>	<b>0.97</b>	<b>0.99</b>	0.96
PQA-net	0.81	0.82	0.84	0.83	0.91	0.92	0.93	0.91	0.89	0.89	0.95	<b>0.96</b>	0.97	0.96
3D-NSS	0.60	0.67	0.85	0.79	0.80	0.84	0.94	0.93	0.90	0.90	0.96	0.93	<b>0.98</b>	0.94
GMS-3DQA	0.83	0.84	0.91	0.92	0.95	0.95	0.95	0.93	0.96	<b>0.97</b>	<b>0.97</b>	0.93	<b>0.98</b>	0.95
MM-PCQA	<b>0.84</b>	0.83	<b>0.92</b>	<b>0.93</b>	0.94	0.96	0.94	<b>0.95</b>	0.94	0.93	0.95	0.94	<b>0.98</b>	0.96
LMM-PCQA	<b>0.89</b>	<b>0.90</b>	<b>0.97</b>	<b>0.97</b>	<b>0.97</b>	<b>0.98</b>	0.95	<b>0.94</b>	<b>0.98</b>	<b>0.98</b>	0.96	0.93	<b>0.98</b>	<b>0.97</b>

**Table 6: Distortion-specific performance results on the WPC database, where DS represents down-sampling, GN represents geometry and color Gaussian noise, G-PCC represents geometry-based compression, and V-PCC represents video-based compression.**

Distortion	DS		GN		G-PCC		V-PCC	
	SR↑	PL↑	SR↑	PL↑	SR↑	PL↑	SR↑	PL↑
MSE-p2po	0.46	0.46	0.67	0.63	0.84	0.72	0.41	0.44
HD-p2po	0.39	0.35	0.69	0.70	0.80	0.70	0.36	0.36
MSE-p2pl	0.41	0.40	0.61	0.47	0.61	0.48	0.42	0.40
HD-p2pl	0.40	0.41	0.62	0.51	0.57	0.59	0.41	0.41
PSNR-yuv	0.23	0.28	0.79	0.88	0.47	0.43	0.44	0.48
PCQM	0.66	0.64	<b>0.89</b>	<b>0.89</b>	0.86	0.77	0.83	0.78
GraphSIM	0.56	0.57	0.79	0.81	0.75	0.74	0.71	0.70
PointSSIM	0.35	0.34	0.83	0.87	<b>0.91</b>	<b>0.95</b>	0.51	0.41
PQA-net	0.61	0.63	0.77	0.78	0.87	0.88	0.76	0.77
3D-NSS	0.55	0.51	0.81	0.83	0.86	0.87	0.49	0.47
GMS-3DQA	0.72	0.73	0.87	0.88	0.89	0.89	0.86	<b>0.88</b>
MM-PCQA	<b>0.75</b>	<b>0.74</b>	0.88	0.87	0.88	0.89	<b>0.89</b>	0.86
LMM-PCQA	<b>0.76</b>	<b>0.79</b>	<b>0.91</b>	<b>0.91</b>	<b>0.96</b>	<b>0.96</b>	<b>0.90</b>	<b>0.89</b>

#### 4.8 Distortion-specific Evaluation

To verify the effectiveness of the proposed LMM-PCQA on different kinds of distortions, we carry out the distortion-specific evaluation experiment in this section. The performance comparison is shown in Table 5 (The ‘nan’ values for MSE-p2po, HD-p2po, MSE-p2pl, and HD-p2pl arise because these methods only analyze the geometric differences between the reference and distorted point clouds, thereby failing to account for color distortions.) and Table 6, from which we can obtain several findings: 1) The proposed LMM-PCQA achieves the best performance on 4 of 7 distortion types and all distortion types on the SJTU-PCQA database and the WPC database respectively, which suggests that LMM-PCQA is effective at dealing with various kinds of distortions. 2) Although LMM-PCQA does not

achieve the top performance on DS+CN, GGN, and CN+GGN distortions within the SJTU-PCQA database, the performance gap to the best is minimal, with a difference of no more than 0.02 in terms of SRCC values. This suggests that LMM-PCQA remains in the top tier. 3) Upon examining Table 6 more closely, it is evident that all PCQA methods exhibit substantial declines in performance with the ‘DS’ distortion. The underlying reason for this trend is the simplicity of the point cloud reference models used in the WPC database, leading to a reduced sensitivity to downsampling distortion.

## 5 CONCLUSION

In conclusion, this paper pioneers the integration of LMMs with PCQA, unveiling the untapped potential of LMMs in this domain. Our research successfully demonstrates the feasibility of adapting LMMs for PCQA through text supervision, enhancing their capability to evaluate 3D visual quality from 2D projections. We also propose to capture multi-scale structural features to offer a more holistic view of the point cloud quality. By combining LMM evaluation results and structural features, our approach significantly improves the accuracy of PCQA. We carry out thorough experiments to demonstrate the effectiveness of the proposed LMM-PCQA, as well as its robustness and ability to generalize across various distortion types and diverse point cloud content. The encouraging results not only validate our methodology but also lay the groundwork for future research. We hope our work will serve as a stepping stone for further exploration into the synergy between LMMs and PCQA, driving forward innovations in 3D visual quality analysis.

## ACKNOWLEDGMENTS

This work was supported in part by the National Natural Science Foundation of China (623B2073, 62101326, 62225112, 62301316), the Foundation of Key Laboratory of Media Audio & Video (Communication University of China), Ministry of Education, China (JGK-FKT2302), the China Postdoctoral Science Foundation under Grants 2023TQ0212 and 2023M742298, and the Postdoctoral Fellowship Program of CPSF under Grant GZC20231618.



## REFERENCES

- [1] 2000. Recommendation 500-10: Methodology for the subjective assessment of the quality of television pictures. ITU-R Rec. BT.500.
- [2] Evangelos Alexiou et al. 2019. Exploiting user interactivity in quality assessment of point cloud imaging. In *QoMEX*. IEEE.
- [3] Evangelos Alexiou and Touradj Ebrahimi. 2020. Towards a point cloud structural similarity metric. In *International Conference on Multimedia and Expo Workshop*. 1–6.
- [4] Evangelos Alexiou, Xuemei Zhou, Irene Viola, and Pablo Cesar. 2021. PointPCA: Point cloud objective quality assessment using PCA-based descriptors. *arXiv preprint arXiv:2111.12663* (2021).
- [5] Xiongli Chai, Feng Shao, Baoyang Mu, Hangwei Chen, Qiuping Jiang, and Yousung Ho. 2024. Plain-PCQA: No-Reference Point Cloud Quality Assessment by Analysis of Plain Visual and Geometrical Components. *IEEE Transactions on Circuits and Systems for Video Technology* (2024).
- [6] Qi Chen, Lin Sun, Zhixin Wang, Kui Jia, and Alan Yuille. 2020. Object as hotspots: An anchor-free 3d object detection approach via firing of hotspots. In *European Conference on Computer Vision*. 68–84.
- [7] Yixiong Chen. 2023. X-IQE: eXplainable Image Quality Evaluation for Text-to-Image Generation with Visual Large Language Models. *arXiv preprint arXiv:2305.10843* (2023).
- [8] Mingmei Cheng, Le Hui, Jin Xie, and Jian Yang. 2021. SSPC-Net: Semi-supervised semantic 3D point cloud segmentation network. In *AAAI*.
- [9] Aladine Chetouani, Maurice Quach, Giuseppe Valenzise, and Frédéric Dufaux. 2021. Deep learning-based quality assessment of 3d point clouds without reference. In *International Conference on Multimedia and Expo Workshop*. 1–6.
- [10] Yaodong Cui, Ren Chen, Wenbo Chu, Long Chen, Daxin Tian, Ying Li, and Dongpu Cao. 2021. Deep learning for image and point cloud fusion in autonomous driving: A review. *IEEE Transactions on Intelligent Transportation Systems* 23, 2 (2021), 722–739.
- [11] Tingyu Fan, Linyao Gao, Yiling Xu, Zhu Li, and Dong Wang. 2022. D-DPCC: Deep Dynamic Point Cloud Compression via 3D Motion Prediction. *International Joint Conference on Artificial Intelligence* (2022).
- [12] Yu Fan, Zicheng Zhang, Wei Sun, Xiongkuo Min, Ning Liu, Quan Zhou, Jun He, Qiyuan Wang, and Guangtao Zhai. 2022. A no-reference quality assessment metric for point cloud based on captured video sequences. In *IEEE MMSP*. IEEE, 1–5.
- [13] Eleonora Grilli, Fabio Menna, and Fabio Remondino. 2017. A review of point clouds segmentation and classification algorithms. *The International Archives of Photogrammetry, Remote Sensing and Spatial Information Sciences* 42 (2017), 339.
- [14] Ke Gu, Dacheng Tao, Jun-Fei Qiao, and Weisi Lin. 2017. Learning a no-reference quality assessment model of enhanced images with big data. *IEEE Transactions on Neural Networks and Learning Systems* 29, 4 (2017), 1301–1313.
- [15] Shuai Gu, Junhui Hou, Huanqiang Zeng, Hui Yuan, and Kai-Kuang Ma. 2019. 3D point cloud attribute compression using geometry-guided sparse representation. *IEEE Transactions on Image Processing* 29 (2019), 796–808.
- [16] Zhipeng Huang, Zhizheng Zhang, Yiting Lu, Zheng-Jun Zha, Zhibo Chen, and Baining Guo. 2024. VisualCritic: Making LMMs Perceive Visual Quality Like Humans. *arXiv:2403.12806* [cs.CV]
- [17] Jason Ku, Melissa Mozifian, Jungwook Lee, Ali Harakeh, and Steven L Waslander. 2018. Joint 3d proposal generation and object detection from view aggregation. In *IEEE/RISJ International Conference on Intelligent Robots and Systems*. 1–8.
- [18] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual Instruction Tuning.
- [19] Qi Liu, Honglei Su, Zhengfang Duanmu, Wentao Liu, and Zhou Wang. 2022. Perceptual Quality Assessment of Colored 3D Point Clouds. *IEEE Transactions on Visualization and Computer Graphics* (2022).
- [20] Qi Liu, Hui Yuan, Raouf Hamzaoui, Honglei Su, Junhui Hou, and Huan Yang. 2021. Reduced reference perceptual quality model with application to rate control for video-based point cloud compression. *IEEE Transactions on Image Processing* 30 (2021), 6623–6636.
- [21] Qi Liu, Hui Yuan, Junhui Hou, Raouf Hamzaoui, and Honglei Su. 2020. Model-based joint bit allocation between geometry and color for video-based 3D point cloud compression. *IEEE Transactions on Multimedia* 23 (2020), 3278–3291.
- [22] Qi Liu, Hui Yuan, Honglei Su, Hao Liu, Yu Wang, Huan Yang, and Junhui Hou. 2021. PQA-Net: Deep No Reference Point Cloud Quality Assessment via Multi-View Projection. *IEEE Transactions on Circuits and Systems for Video Technology* 31, 12 (2021), 4645–4660.
- [23] Tsung-Jung Liu, Kuan-Hsien Liu, Joe Yuchieh Lin, Weisi Lin, and C-C Jay Kuo. 2015. A paraboost method to image quality assessment. *IEEE Transactions on Neural Networks and Learning Systems* 28, 1 (2015), 107–121.
- [24] Weiquan Liu, Hanyun Guo, Weini Zhang, Yu Zang, Cheng Wang, and Jonathan Li. 2022. TopoSeg: Topology-aware Segmentation for Point Clouds. *International Joint Conference on Artificial Intelligence* (2022).
- [25] Yipeng Liu, Qi Yang, Yiling Xu, and Le Yang. 2022. Point Cloud Quality Assessment: Dataset Construction and Learning-based No-Reference Metric. *ACM Transactions on Multimedia Computing, Communications, and Applications* (2022).
- [26] Rufael Mekuria, Kees Blom, and Pablo Cesar. 2016. Design, implementation, and evaluation of a point cloud codec for tele-immersive video. *IEEE Transactions on Circuits and Systems for Video Technology* 27, 4 (2016), 828–842.
- [27] R Mekuria, Z Li, C Tulvan, and P Chou. 2016. Evaluation criteria for point cloud compression. *ISO/IEC MPEG 16332* (2016).
- [28] Gabriel Meynet, Yana Nehmé, Julie Digne, and Guillaume Lavoué. 2020. PCQM: A full-reference quality metric for colored 3D point clouds. In *International Workshop on Quality of Multimedia*. 1–6.
- [29] Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. 2012. No-reference image quality assessment in the spatial domain. *IEEE Transactions on Image Processing* 21, 12 (2012), 4695–4708.
- [30] Anish Mittal, Rajiv Soundararajan, and Alan C Bovik. 2012. Making a “completely blind” image quality analyzer. *IEEE Signal Processing Letters* 20, 3 (2012), 209–212.
- [31] Youngmin Park, Vincent Lepetit, and Woonatck Woo. 2008. Multiple 3d object tracking for augmented reality. In *IEEE/ACM International Symposium on Mixed and Augmented Reality*. 117–120.
- [32] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision.
- [33] Alec Radford, Jeff Wu, Rewon Child, D. Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- [34] Dong Tian, Hideaki Ochimizu, Chen Feng, Robert Cohen, and Anthony Vetro. 2017. Geometric distortion metrics for point cloud compression. In *IEEE International Conference on Image Processing*. 3460–3464.
- [35] Eric M Torlig, Evangelos Alexiou, Tiago A Fonseca, Ricardo L de Queiroz, and Touradj Ebrahimi. 2018. A novel methodology for quality assessment of voxelized point clouds. In *Applications of Digital Image Processing XLI*, Vol. 10752. 174–190.
- [36] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yas-mine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucu-rull, David Esiohu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open Foundation and Fine-Tuned Chat Models. *arXiv:2307.09288* [cs.CL]
- [37] Sourabh Vora, Alex H Lang, Bassam Helou, and Oscar Beijbom. 2020. Pointpainting: Sequential fusion for 3d object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*. 4604–4612.
- [38] Songtao Wang, Xiaoqi Wang, Hao Gao, and Jian Xiong. 2023. Non-Local Geometry and Color Gradient Aggregation Graph Model for No-Reference Point Cloud Quality Assessment. In *Proceedings of the 31st ACM International Conference on Multimedia*. 6803–6810.
- [39] Zhixin Wang and Kui Jia. 2019. Frustum convnet: Sliding frustums to aggregate local point-wise features for amodal 3d object detection. In *IEEE/RISJ International Conference on Intelligent Robots and Systems*. 1742–1749.
- [40] Haoning Wu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Annan Wang, Chunyi Li, Wenxiu Sun, Qiong Yan, Guangtao Zhai, and Weisi Lin. 2024. Q-Bench: A Benchmark for General-Purpose Foundation Models on Low-level Vision. *ICLR* (2024).
- [41] Haoning Wu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Annan Wang, Kaixin Xu, Chunyi Li, Jingwen Hou, Guangtao Zhai, et al. 2024. Q-instruct: Improving low-level visual abilities for multi-modality foundation models. *CVPR* (2024).
- [42] Haoning Wu, Zicheng Zhang, Weixia Zhang, Chaofeng Chen, Liang Liao, Chunyi Li, Yixuan Gao, Annan Wang, Erli Zhang, Wenxiu Sun, et al. 2023. Q-align: Teaching llms for visual scoring via discrete text-defined levels. *arXiv preprint arXiv:2312.17090* (2023).
- [43] Haoning Wu, Hanwei Zhu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Chunyi Li, Annan Wang, Wenxiu Sun, Qiong Yan, Xiaohong Liu, Guangtao Zhai, Shiqi Wang, and Weisi Lin. 2024. Towards Open-ended Visual Quality Comparison. *arXiv preprint arXiv:2402.16641* (2024).
- [44] Liang Xie, Chao Xiang, Zhengxu Yu, Guodong Xu, Zheng Yang, Deng Cai, and Xiaofei He. 2020. PI-RCNN: An efficient multi-sensor 3D object detector with point-based attentive cont-conv fusion module. In *AAAI*, Vol. 34. 12460–12467.
- [45] Wuyuan Xie, Kaimin Wang, Yakun Ju, and Miaohui Wang. 2023. pmbqa: Projection-based blind point cloud quality assessment via multimodal learning. In *Proceedings of the 31st ACM International Conference on Multimedia*. 3250–3258.
- [46] Qi Yang, Hao Chen, Zhan Ma, Yiling Xu, Rongjun Tang, and Jun Sun. 2020. Predicting the perceptual quality of point cloud: A 3d-to-2d projection-based exploration. *IEEE Transactions on Multimedia* (2020).

- [47] Qi Yang, Yipeng Liu, Siheng Chen, Yiling Xu, and Jun Sun. 2022. No-Reference Point Cloud Quality Assessment via Domain Adaptation. In *Proceedings of the IEEE/CVF international conference on computer vision*. 21179–21188.
- [48] Qi Yang, Zhan Ma, Yiling Xu, Zhu Li, and Jun Sun. 2020. Inferring point cloud quality via graph similarity. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2020).
- [49] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, Chaoya Jiang, Chenliang Li, Yuanhong Xu, Hehong Chen, Junfeng Tian, Qian Qi, Ji Zhang, and Fei Huang. 2023. mPLUG-Owl: Modularization Empowers Large Language Models with Multimodality. arXiv:2304.14178 [cs.CL]
- [50] Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. 2023. mPLUG-Owl2: Revolutionizing Multi-modal Large Language Model with Modality Collaboration. arXiv:2311.04257 [cs.CL]
- [51] Jin Hyeok Yoo, Yecheol Kim, Jisong Kim, and Jun Won Choi. 2020. 3d-cvf: Generating joint camera and lidar features using cross-view spatial feature fusion for 3d object detection. In *European Conference on Computer Vision*. 720–736.
- [52] Zhiyuan You, Zheyuan Li, Jinjin Gu, Zhenfei Yin, Tianfan Xue, and Chao Dong. 2023. Depicting Beyond Scores: Advancing Image Quality Assessment through Multi-modal Language Models. arXiv:2312.08962 [cs.CV]
- [53] Chaofan Zhang, Ziqing Huang, Shiguang Liu, and Jian Xiao. 2022. Dual-Channel Multi-Task CNN for No-Reference Screen Content Image Quality Assessment. *IEEE Transactions on Circuits and Systems for Video Technology* 32, 8 (2022), 5011–5025.
- [54] Chaofan Zhang and Shiguang Liu. 2022. No-reference omnidirectional image quality assessment based on joint network. In *ACM International Conference on Multimedia*. 943–951.
- [55] Lin Zhang, Lei Zhang, and Alan C Bovik. 2015. A feature-enriched completely blind image quality evaluator. *IEEE Transactions on Image Processing* 24, 8 (2015), 2579–2591.
- [56] Wei Zhang, Ali Borji, Zhou Wang, Patrick Le Callet, and Hantao Liu. 2015. The application of visual saliency models in objective image quality assessment: A statistical evaluation. *IEEE Transactions on Neural Networks and Learning Systems* 27, 6 (2015), 1266–1278.
- [57] Yujie Zhang, Qi Yang, and Yiling Xu. 2021. MS-GraphSIM: Inferring point cloud quality via multiscale graph similarity. In *Proceedings of the 29th ACM International Conference on Multimedia*. 1230–1238.
- [58] Zicheng Zhang, Wei Sun, Xiongkuo Min, Tao Wang, Wei Lu, and Guangtao Zhai. 2022. No-reference quality assessment for 3d colored point cloud and mesh models. *IEEE Transactions on Circuits and Systems for Video Technology* (2022).
- [59] Zicheng Zhang, Wei Sun, Xiongkuo Min, Wei Wu, Ying Chen, and Guangtao Zhai. 2022. Treating Point Cloud as Moving Camera Videos: A No-Reference Quality Assessment Metric. *arXiv preprint arXiv:2208.14085* (2022).
- [60] Zicheng Zhang, Wei Sun, Xiongkuo Min, Quan Zhou, Jun He, Qiyuan Wang, and Guangtao Zhai. 2023. MM-PCQA: Multi-modal learning for no-reference point cloud quality assessment. *International Joint Conference on Artificial Intelligence* (2023).
- [61] Zicheng Zhang, Wei Sun, Xiongkuo Min, Wenhan Zhu, Tao Wang, Wei Lu, and Guangtao Zhai. 2021. A No-Reference Evaluation Metric for Low-Light Image Enhancement. In *IEEE International Conference on Multimedia and Expo*.
- [62] Zicheng Zhang, Wei Sun, Houning Wu, Yingjie Zhou, Chunyi Li, Xiongkuo Min, Guangtao Zhai, and Weisi Lin. 2023. GMS-3DQA: Projection-based Grid Minipatch Sampling for 3D Model Quality Assessment. *arXiv preprint arXiv:2306.05658* (2023).
- [63] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. 2018. Open3D: A modern library for 3D data processing. *arXiv preprint arXiv:1801.09847* (2018).
- [64] Wei Zhou, Qi Yang, Qiuping Jiang, Guangtao Zhai, and Weisi Lin. 2022. Blind Quality Assessment of 3D Dense Point Clouds with Structure Guided Resampling. *arXiv preprint arXiv:2208.14603* (2022).
- [65] Xuemei Zhou, Evangelos Alexiou, Irene Viola, and Pablo Cesar. 2023. PointPCA+: Extending PointPCA objective quality assessment metric. In *2023 IEEE International Conference on Image Processing Challenges and Workshops (ICIPCW)*. IEEE, 1–5.