# Ultrasound-QBench: Can LLMs Aid in Quality Assessment of Ultrasound Imaging?

Hongyi Miao[1], Junjia Liu[2]†, Yankun Cao[1], Yingjie Zhou[2], Yanwei Jiang[2], Zhi Liu[1]†, Guangtao Zhai[2]
[1]Shandong University, [2]Shanghai Jiao Tong University
†Corresponding author

*Abstract*—With the dramatic upsurge in the volume of ultrasound examinations, low-quality ultrasound imaging has gradually increased due to variations in operator proficiency and imaging circumstances, imposing a severe burden on diagnosis accuracy and even entailing the risk of restarting the diagnosis in critical cases. To assist clinicians in selecting high-quality ultrasound images and ensuring accurate diagnoses, we introduce Ultrasound-QBench, a comprehensive benchmark that systematically evaluates multimodal large language models (MLLMs) on quality assessment tasks of ultrasound images. Ultrasound-QBench establishes two datasets collected from diverse sources: IVUSQA, consisting of 7,709 images, and CardiacUltraQA, containing 3,863 images. These images encompassing common ultrasound imaging artifacts are annotated by professional ultrasound experts and classified into three quality levels: high, medium, and low. To better evaluate MLLMs, we decompose the quality assessment task into three dimensionalities: qualitative classification, quantitative scoring, and comparative assessment. The evaluation of 7 open-source MLLMs as well as 1 proprietary MLLMs demonstrates that MLLMs possess preliminary capabilities for low-level visual tasks in ultrasound image quality classification. We hope this benchmark will inspire the research community to delve deeper into uncovering and enhancing the untapped potential of MLLMs for medical imaging tasks.

*Index Terms*—Multimodal Large Language Model (MLLM), Quality Assessment, Ultrasound Image

## I. INTRODUCTION

Ultrasound imaging represents a medical imaging technology that is prevalently utilized in clinical practice. The advantages of ultrasound imaging such as low cost, ease of operation, radiation-free nature, and the ability to provide real-time imaging, render it ideal for visualizing soft tissues [1]. Currently, ultrasound imaging has been extensively employed for diagnosing abdominal, cardiac, vascular, and musculoskeletal diseases, as well as for prenatal examinations. Nevertheless, with the increasing volume of daily ultrasound examinations, the variability in image quality has emerged as a considerable challenge that affects diagnostic accuracy, data management, and healthcare efficiency [2]. Low-quality images not only diminish diagnostic precision but also lead to repeat scans, thereby increasing healthcare costs and resource wastage [3].

In order to evaluate image quality, quality assessment (QA) has been extensively investigated in the domain of natural images. In traditional quality assessment methods, image features are manually extracted and statistically analyzed, and subsequently compared with reference images for full-reference methods such as Structural Similarity Index (SSIM) [4] or directly estimated for no-reference methods such as BRISQUE [5] and NIQE [6]. In contrast to natural images, quality assessment of medical images involves diagnostic relevance and the handling of artifacts specific to medical images, such as noise and blurring [7], rendering traditional methods less effective. Specifically, ultrasound images often incorporate distinctive artifacts resulting from operator variability and diverse imaging conditions [8], [9] such as multiple reflections, multiple internal reverberations, and refractive shadow. Consequently, merely reutilizing the image features employed in natural image quality assessment without taking into account the specialized characteristics of ultrasound imaging is unable to address these challenges [10]. To further enhance the generalization of traditional methods, researchers explore the capabilities of neural networks in quality assessment. These methods convert quality assessment tasks into an end-to-end classification or regression problem, and substitute traditional hand-crafted feature extraction with learning-based feature representations. However, the learning process relies on a considerable amount of labeled data, which is costly and scarce in the context of medical imaging, and the assessment performance may deteriorate in the presence of unlabeled noise and artifacts.

Based on the aforementioned analysis, an ideal quality assessment method of ultrasound images requires to be both specialized and generalizable. The former necessitates that the quality assessment method fully utilize the characteristics of ultrasound imaging, while the latter requires that the quality assessment method be capable of making accurate judgments for different types of artifacts. Recent progress in Multimodal Large Language Models (MLLMs), such as LLaVA [11], Qwen2-vl [12], and mPLUG-owl3 [13], holds great promise in medical image quality assessment due to their zero-shot inference capacity and cross-domain expertise. Unlike traditional neural networks that require a large amount of labeled data for training, MLLMs can leverage generative pre-trained models and given prompts to conduct inference, demonstrating excellent performance with limited labeled data or even in zero-shot scenarios [14]–[17]. This suggests the potential for MLLMs to accurately assess ultrasound image quality, even in the presence of unobserved distortion types. Moreover, the pre-trained models have learned from an extensive amount of labeled data across diverse fields, presenting exceptional cross-domain expertise. Based on these advantages, MLLMs are both specialized in cross-domains and generalizable to
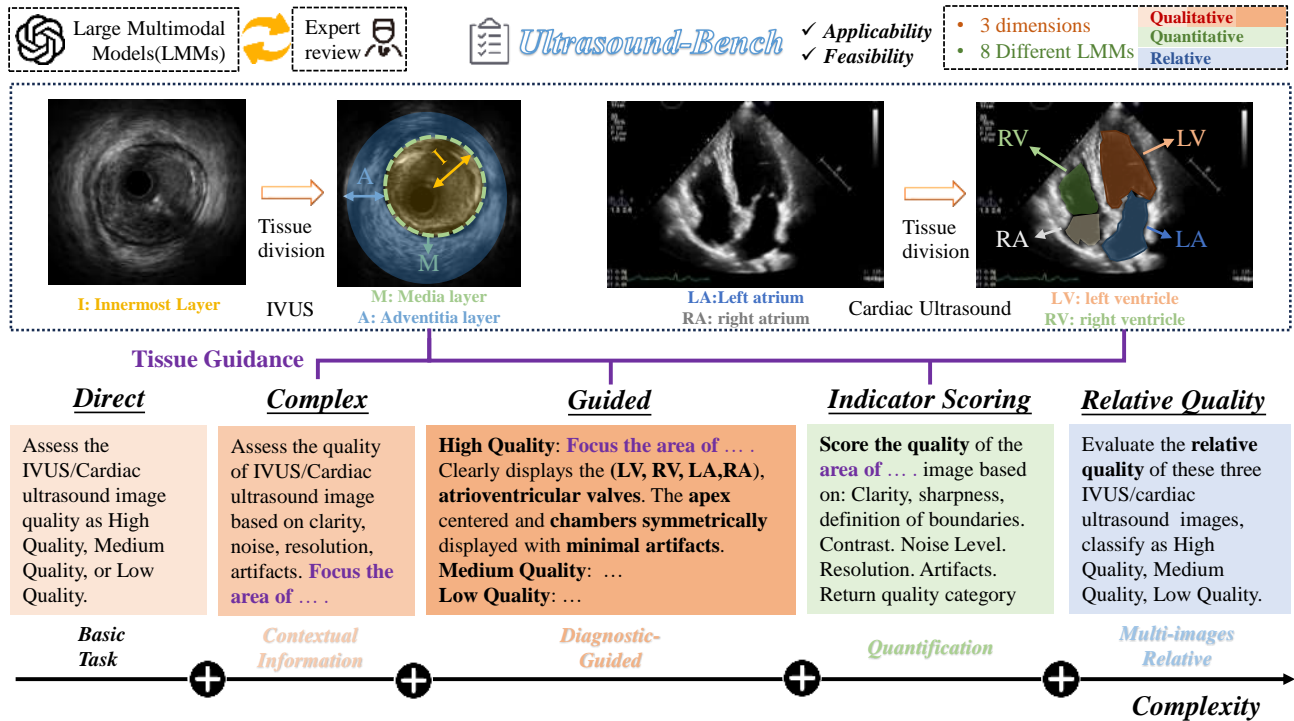
Fig. 1. In the proposed Ultrasound-QBench, we established the first benchmark on MLLM capabilities On ultrasound images, qualitative quality assessment, quantitative evaluation, and relative quality are included.

unencountered samples, making them particularly appropriate for ultrasound image quality assessment.

Inspired by the quality evaluation experiments on natural images using MLLMs [18]–[20], we propose a new benchmark framework, **Ultrasound-QBench**, to evaluate MLLMs in ultrasound image quality assessment. In summary, Ultrasound-QBench investigates the generalization capacity of MLLMs to ultrasound image quality assessment without task-specific fine-tuning, bridging the gap between natural image QA and medical QA. Fig. 1 present the overall diagram of Ultrasound-QBench framework. We evaluate 7 open-source MLLMs as well as 1 proprietary MLLM from three dimensionalities:

- **Qualitative Classification:** MLLMs are required to classify the ultrasound images into three quality degrees: low, medium, and high. Prompts spanning from rough to diagnostic level are utilized to assist MLLMs in making judgments. As depicted in Fig. 1, the prompts are categorized into three complexity levels: basic question, contextual information, and diagnostic-guided information. For the second and third tasks, the tissues information is included in the prompts.
- **Quantitative Scoring:** MLLMs are required to provide quantitative scores in multiple quality-related indicators such as clarity, contrast, noise level, detail resolution, uniformity, and presence of artifact. As demonstrated in Fig. 1, the prompts contain the tissues information, guiding the MLLMs to focus on the region of interest.
- **Comparative Assessment:** MLLMs are required to con-

duct a comparison of the relative quality among multiple ultrasound images. This task involves evaluating the ability of MLLMs to comprehend the relative quality changes between diverse images.

Ultrasound-QBench establishes two datasets collected from diverse sources: IVUSQA, comprising 7,709 images, and CardiacUltraQA, containing 3,863 images. These images encompass common ultrasound imaging artifacts and noises, e.g. multiple reflections, multiple internal reverberations, and refractive shadow, and are annotated by professional ultrasound experts and classified into three quality levels: high, medium, and low. By addressing the limitations of traditional quality assessment methods and leveraging the strengths of MLLMs, this work lays the foundation for advancing automated quality assessment of ultrasound images. Our findings provide valuable insights into the potential of MLLMs for optimizing diagnostic workflows in clinical practice.

## II. DATASET

### A. Overview

To evaluate the performance of multimodal large language models (MLLMs) in ultrasound image quality assessment, we establish two ultrasound image datasets that represent two typical real-world clinical scenarios: IVUSQA, concentrating on intravascular ultrasound (IVUS) images, and CardiacUltraQA, covering cardiac ultrasound images. These two datasets are curated to represent diverse imaging conditions and feature

high-quality expert annotations, providing a robust framework for zero-shot testing.
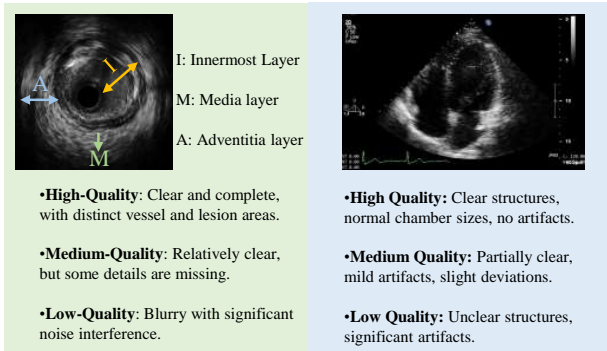


Fig. 2. IVUSQA and CardiacUltraQA Dataset and Assessment Standard.



Fig. 3. Quality Distribution combined for IVUSQA and CardiacUltraQA Dataset.

## B. Dataset Composition

**IVUSQA Dataset:** This dataset contains a total of 7,709 intravascular ultrasound (IVUS) images, which are used to assess vascular structures. IVUS imaging is characterized by a circular view, centered on the vessel wall, and is mainly used to detect plaques and other abnormalities in arteries. The three primary components of an IVUS image include: **(1) Innermost Layer:** this layer comprises the intima, atheroma, and internal elastic membrane; **(2) Media Layer:** composed of smooth muscle cells that do not reflect ultrasound, appearing as dark areas in the image; **(3) Adventitia Layer:** this layer consists of collagen, which reflects a significant amount of ultrasound, presenting as white in IVUS images.

**CardiacUltraQA Dataset:** The CardiacUltraQA dataset consists of 3,863 the Apical Four-Chamber View images. This view displays the four main chambers of the heart (left atrium, right atrium, left ventricle, right ventricle) along with associated structures such as the heart valves and surrounding tissues. These images are crucial for assessing heart health, particularly for observing the symmetry of the heart chambers and normal anatomical features.

The quality of images in these two datasets is classified into three categories: high, medium, and low, based on the standard illustrated in Fig. 2. All images are annotated by a team of certified ultrasound experts using a unified and standardized subjective quality assessment framework to ensure consistency and clinical relevance. Fig. 3 depicts the subjective quality distribution of these two datasets.

## C. Advantages of the Dataset

- **Clinical Significance:** Both IVUSQA and CardiacUltraQA focus on critical diagnostic areas, including vascular health and cardiac function, ensuring the applicability of the evaluation results in clinical practice.
- **Data Diversity:** The images in the dataset originate from various operators and imaging conditions, thereby incorporating a wide range of noise and artifacts commonly encountered in ultrasound imaging.
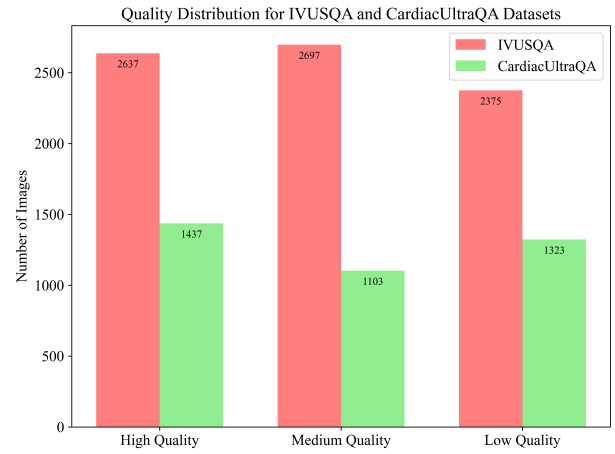
- **Expert Annotations:** The quality label attached to each image is obtained by multiple ultrasound physicians and medical doctoral students according to a standardized and unified subjective assessment process, ensuring high reliability in clinical practice.

## III. EVALUATION WORKFLOW

### A. Prompt Design

To systematically evaluate the performance of Multimodal Large Language Models (MLLMs) in ultrasound image quality assessment, we decompose the evaluating task into three sub-tasks: (1) qualitative classification, (2) quantitative scoring, and (3) comparative assessment.

*a) **Qualitative Classification**:* In this sub-task, MLLMs are required to classify ultrasound images into three quality categories: low, medium, and high. As illustrated in the three orange boxes of Fig. 1, prompts of varying complexity are provided to assist MLLMs in assessing image quality: basic questions, contextual details, and diagnostic guidance. For the second and third tasks, additional tissue information is integrated into the prompts to help MLLMs focus on critical regions. These three types of prompts are described in detail as follows:

- **basic questions:** the model's baseline performance in image quality assessment without additional context or guidance is evaluated.
- **contextual details:** additional contextual information, such as the anatomical region of the ultrasound image or the clinical use case, is provided to the model. The impact of more detailed contextual descriptions on the model's ability to enhance classification accuracy and its understanding of task-specific information is evaluated.
- **diagnostic guidance:** specific diagnostic criteria is introduced to guide the classification process.
  These prompts evaluate the model's ability to follow detailed clinical instructions and make classification decisions based on predefined diagnostic criteria.
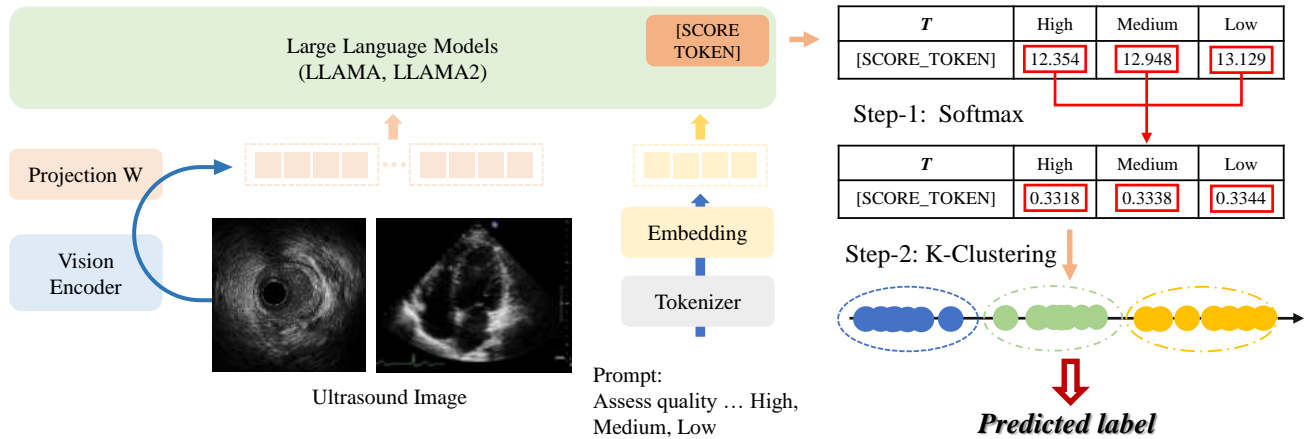
Fig. 4. The proposed softmax-based quality assessment strategy for MLLMs improves upon existing methods by extracting logits for the 'high quality,' 'medium quality,' and 'low quality' categories, rather than directly decoding tokens from the [SCORE TOKEN] position. The strategy predicts labels through a weighted summation and pooling of these logits, followed by a weighted clustering to obtain the final quality rating.

*b) Quantitative Scoring:* This sub-task requires the model to provide an overall quality score based on multiple quality indicators such as clarity, noise level, resolution, thereby evaluating the model's ability to perform multi-dimensional reasoning. The **green** box of Fig. 1 presents a typical template of this kind of prompt.

*c) Comparative Assessment:* This sub-task challenges the model to compare multiple images and identify their relative quality differences. The **blue** box of Fig. 1 presents a typical template of this kind of prompt.

### B. Experimental Workflow

The experimental workflow is organized into three stages: model selection, execution, and evaluation.

*a) Model Selection Stage:* In this stage, we select 7 open-source MLLMs as well as 1 proprietary MLLM to evaluate their capacities of quality assessment for IVUSQA and CardiacUltraQA. These models include LLaVA-v1.5-7b, LLaVA-v1.5-13b, InternLM-XComposer2-VL-7b [21], DeepSeek [22], LLaVA-Med, Qwen2-VL, mPLUG-Owl3, and GPT4o (proprietary). The five kinds of prompts corresponding to the three sub-tasks are predefined to ensure consistency across the experiments. Optimal settings for each model are selected to ensure reliable inference results.

*b) Execution Stage:* In this stage, ultrasound images are processed in batches of 16 on an NVIDIA RTX 4090 GPU with 24GB of memory. For each image-prompt pair, the evaluated model generates a response, which is then parsed to extract the score tokens. To improve the assessment accuracy and prevent predictions from being biased toward extreme outcomes, we propose a new evaluation strategy that combines Softmax and k-Clustering, as illustrated in Fig. 4. The strategy consists of two steps:

- **Step-1 Softmax Strategy:** The model outputs score_tokens for each quality level based on the prompt. By applying the Softmax operation, the raw logits are transformed into comparable probability

values, reducing the impact of extreme values on the model's predictions and preventing the model from excessively favoring certain classes.

- **Step-2 Clustering:** The clustering method is combined with a weighted summation approach to optimize the predictions output by Step-1. By grouping similar prediction results together, clustering effectively prevents the model from excessively favoring any one class, ensuring balanced prediction results. Clustering can automatically adjust the class boundaries based on the distribution of the prediction results, reducing the risk of model bias due to extreme values. This helps to ensure a more balanced distribution of samples across each class and improves the accuracy and stability of the model's classification.

*c) Evaluation Stage:* The evaluation stage computes the key metric of classification accuracy and assesses prompt sensitivity to evaluate the models' performance across different strategies. The analysis includes a comparison of model performance on the two datasets to assess the models' generalizability across different types of ultrasound images. Prompt sensitivity is evaluated by examining how the models' performance varies with different prompt strategies.

## IV. EXPERIMENTAL RESULTS

This section presents and analyze the performance of selected 8 MLLMs on ultrasound image quality assessment. The evaluating results of the selected models on three tasks are presented in Table I.

### A. Original Evaluation without Softmax and Clustering

We first analyze the original performance of MLLMs without Softmax and k-Clustering.

*a) Qualitative Classification:* As shown in Table I, the original capabilities of the selected MLLMs to qualitatively assess ultrasound image quality are disappointing. For basic questions, all of these models exhibit a poor accuracy of around 30% on both IVUSQA and CardiacUltraQA. Through

TABLE I
ACCURACY FOR ULTRASOUND IMAGE QUALITY ASSESSMENT ACROSS TWO DATASETS WITH FIVE PROMPTS. RED INDICATES THE BEST RESULT IN EACH TASK, AND UNDERLINED VALUES REPRESENT THE MOST SIGNIFICANT IMPROVEMENT WITH SOFTMAX + CLUSTERING.

| MLLM | IVUSQA | | | | | CardiacUltraQA | | | | |
| | Qualitative | | | Quantitative | Comparative | Qualitative | | | Quantitative | Comparative |
| | Basic | Contextual | Diagnostic | | | Basic | Contextual | Diagnostic | | |
|---|---|---|---|---|---|---|---|---|---|---|
| LLaVA-v1.5-7B | **34.27%**/46.52% | 34.23%/46.63% | 34.21%/40.36% | 34.31%/<u>48.71%</u> | 94.06% | 37.21%/44.74% | 37.32%/<u>51.51%</u> | 37.26%/**48.44%** | 37.35%/**<u>50.14%</u>** | 97.28% |
| LLaVA-v1.5-13B | 27.14%/47.65% | 34.99%/<u>50.11%</u> | 35.06%45.88% | **34.99%**/<u>49.25%</u> | 97.60% | **40.31%**/46.27% | 31.14%/45.01% | 35.14%/42.58% | **40.37%**/46.33% | 98.45% |
| InternLM-VL | 34.20%/41.47% | 30.81%/41.99% | 34.21%/40.74% | 34.35%/42.86% | 55.64% | 36.86%/44.07% | 37.67%/45.57% | 37.07%/45.46% | 38.15%/44.29% | 80.27% |
| Deepseek | 32.46%/42.38% | 34.27%/43.57% | 32.53%/41.35% | 34.21%/46.37% | 49.64% | 39.71%/42.58% | 37.19%/43.93% | 37.21%/46.74% | 37.20%45.52% | 73.54% |
| Qwen2-VL | 30.90%/41.19% | **36.99%**/42.63% | 30.82%/<u>45.65%</u> | **34.99%**/45.93% | 98.94% | 38.99%/<u>50.01%</u> | **55.94%**/**61.37%** | 34.25%/<u>45.90%</u> | 37.21%/45.18% | 99.35% |
| mPLUG-Owl3 | 31.08%/<u>51.38%</u> | 34.17%/36.53% | **57.88%**/**64.04%** | 34.21%/47.43% | 98.66% | 38.67%/48.49% | 38.78%/46.37% | **44.63%**/45.45% | 38.67%/47.99% | 99.64% |
| LLaVA-Med | 34.21%/46.26% | 34.21%/39.64% | 34.99%/40.03% | 34.21%/41.74% | 96.42% | 38.67%/46.26% | 38.78%/42.63% | 37.21%/43.16% | 37.21%/44.61% | 99.88% |
| GPT-4o (proprietary) | 35.63% | 38.41% | 41.24% | 41.31% | 99.93% | 40.67 | 42.13% | 45.37% | 44.28% | 99.96% |

analyzing these results, we find an obvious bias toward a specific class. For instance, mPLUG-Owl tends to classify all images as high quality. For complex contextual information, Qwen2-VL achieves 55.94% accuracy on CardiacUltraQA, while other models still perform poorly. After introducing the diagnostic-level standard into prompts, only mPLUG-Owl3 can improve its accuracy to 57.88%, with the accuracy of medium-quality and low-quality images reaching approximately 50%. These results indicate that existing MLLMs struggle to understand the classification standard of ultrasound image quality in the same way as humans.

*b) Quantitative Scoring:* As shown in Table I, all of these models perform poorly in predicting accurate quality scores, indicating that existing MLLMs struggle to understand the definitions of quality-related indicators.

*c) Comparative Assessment:* As shown in Table I, all of these models can accurately distinguish the relative quality of the given image sequences, indicating that these models can detect changes in ultrasound image quality.

### B. Evaluation with Softmax and Clustering

We performed an extensive evaluation of the performance of Multimodal Large Language Models (MLLMs) augmented with a Softmax strategy enhanced by k-means clustering, with the experimental results presented in Table I. Prior to the incorporation of this strategy, the output distributions across tasks exhibited significant imbalance. For instance, LLaVA-v1.5-7B tended to predict high-quality classifications predominantly across all tasks, except for the relative evaluation task. Upon integrating k-means clustering with the Softmax approach, the model outputs became more balanced, effectively mitigating the tendency of the model to excessively rely on high-quality predictions. This strategy resulted in an average improvement of approximately 10.11% in accuracy across all evaluation prompts. In particular, mPLUG-Owl3 demonstrated a notable accuracy gain of 20.3% in the basic task. The proposed approach successfully mitigates the adverse effects of imbalanced output distributions, thereby enhancing the overall balance of predictions and increasing the robustness of the model's performance. In Table I, the best performance for each task is marked in **red**, while the most significant

improvements attributed to this strategy are **<u>underlined</u>** for each task.

## V. DISCUSSION

### A. Limitations

Despite the improvements achieved with Softmax + Clustering, the performance of MLLMs in quantitative and qualitative assessment of ultrasound image quality remains unsatisfactory. For qualitative assessment, these models still exhibit limitations in assessment accuracy, particularly in distinguishing between medium and low-quality images. For instance, models such as mPLUG-Owl have a tendency to overclassify images as high quality, which indicates an inadequate ability to accurately assess different quality levels. In terms of quantitative scoring, all models fail to predict accurate quality scores reliably. This highlights the necessity for enhanced domain-specific knowledge and advanced feature extraction techniques to improve overall performance. Furthermore, we believe there are two additional factors that limit the performance of MLLMs on the task of ultrasound image quality assessment:

*a) Insufficient Understanding of Ultrasound Image Features:* Despite prompt guidance, current MLLMs lack a full understanding of ultrasound-specific features such as speckle noise, operator variability, and artifacts. These features complicate accurate classification, particularly for high-quality images. Fine-tuning with high-quality, domain-specific data is necessary for models to better capture these unique characteristics.

*b) Dependence on Prompt Engineering:* MLLMs heavily rely on carefully designed prompts, which limits their adaptability in real-world scenarios where input quality can vary. Ultrasound images exhibit significant differences in resolution, contrast, and artifacts. This reliance reduces model flexibility, causing performance degradation when faced with simpler or less structured prompts. Enhancing model adaptability is crucial for real-world medical image assessments.

### B. Future Research Directions

Future research should focus on the following areas:

*a) Reducing Dependence on Prompt Engineering:* Developing adaptive methods to reduce reliance on specific prompts can enhance model flexibility and applicability in real-world medical scenarios.

*b) Expanding Dataset Diversity:* In order to improve model generalization and performance on low-quality images, we will further expand ultrasound datasets to include diverse imaging modalities, patient groups, and clinical conditions.

*c) Leveraging Domain-Specific Knowledge:* By fine-tuning with high-quality labeled data, we can incorporate ultrasound-specific knowledge into model training. Additionally, exploring techniques such as few-shot learning and self-supervised learning can improve performance in data-limited scenarios, thereby enhancing both qualitative and quantitative assessments.

## VI. CONCLUSION

This paper presents Ultrasound-QBench, a benchmark for evaluating multimodal large language models (MLLMs) in ultrasound image quality assessment using the established IVUSQA and CardiacUltraQA datasets. Eight MLLMs are evaluated across qualitative classification, quantitative scoring, and comparative assessment tasks. While models like mPLUG-Owl and Qwen2 show potential in judging ultrasound image quality like humans, they struggle with accurately distinguishing between image quality levels and comprehending ultrasound-specific features, such as distortions and noise. The Softmax + Clustering method can improve accuracy by 10.11%, but limitations remain in understanding ultrasound image structure.

Future work should focus on reducing dependency on prompt engineering, enhancing dataset diversity to improve generalization, and incorporating domain-specific knowledge. Fine-tuning with high-quality labeled data, along with techniques like few-shot learning and self-supervised learning, will further strengthen MLLMs' performance in ultrasound image quality assessment.

## REFERENCES

[1] Lingyun Wu, Jie-Zhi Cheng, Shengli Li, Baiying Lei, Tianfu Wang, and Dong Ni, "Fuiqa: fetal ultrasound image quality assessment with deep convolutional networks," *IEEE transactions on cybernetics*, vol. 47, no. 5, pp. 1336–1349, 2017.

[2] Anton V Nikolaev, Leon De Jong, Gert Weijers, Vincent Groenhuis, Ritse M Mann, Françoise J Siepel, Bogdan M Maris, Stefano Stramigioli, Hendrik HG Hansen, and Chris L De Korte, "Quantitative evaluation of an automated cone-based breast ultrasound scanner for mri–3d us image fusion," *IEEE transactions on medical imaging*, vol. 40, no. 4, pp. 1229–1239, 2021.

[3] Yuxin Song, Zhaoming Zhong, Baoliang Zhao, Peng Zhang, Qiong Wang, Ziwen Wang, Liang Yao, Faqin Lv, and Ying Hu, "Medical ultrasound image quality assessment for autonomous robotic screening," *IEEE Robotics and Automation Letters*, vol. 7, no. 3, pp. 6290–6296, 2022.

[4] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.

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

[6] Anish Mittal, Rajiv Soundararajan, and Alan C Bovik, "Making a "completely blind" image quality analyzer," *IEEE Signal processing letters*, vol. 20, no. 3, pp. 209–212, 2012.

[7] Rafael Rodrigues, Lucie Lévêque, Jesús Gutiérrez, Houda Jebbari, Meriem Outtas, Lu Zhang, Aladine Chetouani, Shaymaa Al-Juboori, Maria G Martini, and Antonio MG Pinheiro, "Objective quality assessment of medical images and videos: Review and challenges," *Multimedia Tools and Applications*, pp. 1–34, 2024.

[8] Qi Chen, Xiongkuo Min, Huiyu Duan, Yucheng Zhu, and Guangtao Zhai, "Muiqa: Image quality assessment database and algorithm for medical ultrasound images," in *2021 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2021, pp. 2958–2962.

[9] Kohei Ohashi, Yukihiro Nagatani, Makoto Yoshigoe, Kyohei Iwai, Keiko Tsuchiya, Atsunobu Hino, Yukako Kida, Asumi Yamazaki, and Takayuki Ishida, "Applicability evaluation of full-reference image quality assessment methods for computed tomography images," *Journal of Digital Imaging*, vol. 36, no. 6, pp. 2623–2634, 2023.

[10] Jinbao Dong, Shengfeng Liu, Yimei Liao, Huaxuan Wen, Baiying Lei, Shengli Li, and Tianfu Wang, "A generic quality control framework for fetal ultrasound cardiac four-chamber planes," *IEEE journal of biomedical and health informatics*, vol. 24, no. 4, pp. 931–942, 2019.

[11] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee, "Visual instruction tuning," *Advances in neural information processing systems*, vol. 36, 2024.

[12] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al., "Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution," *arXiv preprint arXiv:2409.12191*, 2024.

[13] Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou, "mplug-owl3: Towards long image-sequence understanding in multi-modal large language models," *arXiv preprint arXiv:2408.04840*, 2024.

[14] Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao, "Llava-med: Training a large language-and-vision assistant for biomedicine in one day," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[15] Khaled Saab, Tao Tu, Wei-Hung Weng, Ryutaro Tanno, David Stutz, Ellery Wulczyn, Fan Zhang, Tim Strother, Chunjong Park, Elahe Vedadi, et al., "Capabilities of gemini models in medicine," *arXiv preprint arXiv:2404.18416*, 2024.

[16] Xiangru Tang, Anni Zou, Zhuosheng Zhang, Ziming Li, Yilun Zhao, Xingyao Zhang, Arman Cohan, and Mark Gerstein, "Medagents: Large language models as collaborators for zero-shot medical reasoning," *arXiv preprint arXiv:2311.10537*, 2023.

[17] Zhihong Chen, Shizhe Diao, Benyou Wang, Guanbin Li, and Xiang Wan, "Towards unifying medical vision-and-language pre-training via soft prompts," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 23403–23413.

[18] H. Wu, Z. Zhang, E. Zhang, C. Chen, L. Liao, A. Wang, C. Li, W. Sun, Q. Yan, G. Zhai, and W. Lin, "Q-bench: A benchmark for general-purpose foundation models on low-level vision," in *Proceedings of the International Conference on Learning Representations (ICLR)*, Apr. 2024, Spotlight Presentation.

[19] H. Wu, Z. Zhang, E. Zhang, C. Chen, L. Liao, A. Wang, and et al., "Q-instruct: Improving low-level visual abilities for multi-modality foundation models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2024, pp. 1234–1245.

[20] H. Wu, Z. Zhang, W. Zhang, C. Chen, L. Liao, C. Li, and et al., "Q-align: Teaching lmms for visual scoring via discrete text-defined levels," in *Proceedings of the International Conference on Machine Learning (ICML)*, Vienna, Austria, July 2024, pp. 5678–5689.

[21] X. Dong, P. Zhang, Y. Zang, Y. Cao, B. Wang, L. Ouyang, X. Wei, S. Zhang, H. Duan, M. Cao, W. Zhang, Y. Li, H. Yan, Y. Gao, X. Zhang, W. Li, J. Li, K. Chen, C. He, X. Zhang, Y. Qiao, D. Lin, and J. Wang, "Internlm-xcomposer2: Mastering free-form text-image composition and comprehension in vision-language large model," in *Proceedings of the Computer Vision and Pattern Recognition Conference*, Jan. 2024, pp. 1–18, Virtual.

[22] A. Liu, B. Feng, B. Wang, B. Wang, B. Liu, C. Zhao, C. Dengr, C. Ruan, D. Dai, D. Guo, and et al., "Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model," *arXiv preprint arXiv:2405.04434*, 2024.