

MEDQ-BENCH: EVALUATING AND EXPLORING MEDICAL IMAGE QUALITY ASSESSMENT ABILITIES IN MLLMs

Jiayao Liu^{1*}, Jinjie Wei^{1*}, Wanying Qu¹, Chenglong Ma^{1,2}, Junzhi Ning², Yunheng Li¹, Ying Chen², Xinzhe Luo³, Pengcheng Chen², Xin Gao¹, Ming Hu², Huihui Xu², Xin Wang², Shujian Gao¹, Dingkang Yang¹, Zhongying Deng⁴, Jin Ye², Lihao Liu^{2†}, Junjun He^{2†}, Ningsheng Xu¹

¹Fudan University, ²Shanghai Artificial Intelligence Laboratory, ³Imperial College London,

⁴University of Cambridge

*Equal contribution. †Corresponding author. Project Page: <https://github.com/liujiyao/FDU/MedQBench>.

ABSTRACT

Medical Image Quality Assessment (IQA) serves as the first-mile safety gate for clinical AI, yet existing approaches remain constrained by scalar, score-based metrics and fail to reflect the descriptive, human-like reasoning process central to expert evaluation. To address this gap, we introduce **MedQ-Bench**, a comprehensive benchmark that establishes a **perception–reasoning paradigm** for language-based evaluation of medical image quality with Multi-modal Large Language Models (MLLMs). **MedQ-Bench** defines two complementary tasks: (1) **MedQ-Perception**, which probes low-level perceptual capability via human-curated questions on fundamental visual attributes; and (2) **MedQ-Reasoning**, encompassing both *no-reference* and *comparison reasoning* tasks, aligning model evaluation with human-like reasoning on image quality. The benchmark spans *5 imaging modalities* and *over 40 quality attributes*, totaling *2,600 perceptual queries* and *708 reasoning assessments*, covering diverse image sources including authentic clinical acquisitions, images with simulated degradations via physics-based reconstructions, and AI-generated images. To evaluate reasoning ability, we propose a *multi-dimensional judging protocol* that assesses model outputs along four complementary axes. We further conduct rigorous *human–AI alignment validation* by comparing LLM-based judgement with radiologists. Our evaluation of *14 state-of-the-art MLLMs* demonstrates that models exhibit preliminary but unstable perceptual and reasoning skills, with insufficient accuracy for reliable clinical use. These findings highlight the need for targeted optimization of MLLMs in medical IQA. We hope that MedQ-Bench will catalyze further exploration and unlock the untapped potential of MLLMs for medical image quality evaluation.

1 INTRODUCTION

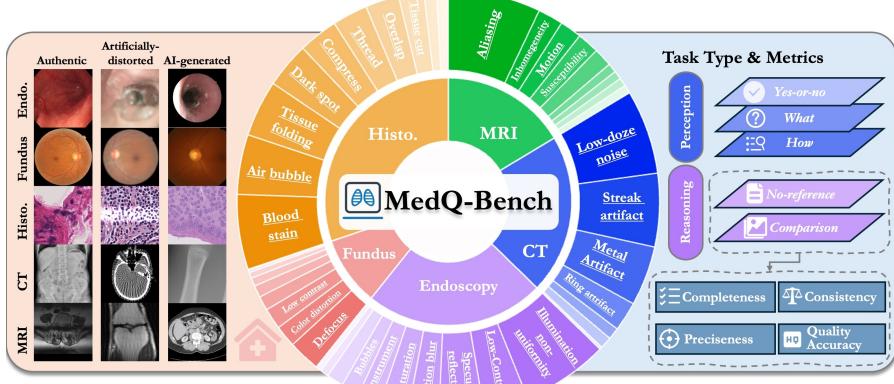


Figure 1: **MedQ-Bench** overview, evaluating MLLMs’ abilities in medical image quality assessment with: (1) Comprehensive coverage: 3,308 samples across 5 modalities with 40+ degradation types. (2) Multi-faceted evaluation: perception-reasoning paradigm.

054 Medical Image Quality Assessment (IQA) determines whether imaging data can be reliably used for
 055 subsequent diagnostic interpretation and clinical decision-making (Lamard et al., 2024). In clinical
 056 practice, multiple visual quality attributes of medical images directly influence diagnostic accuracy
 057 and patient safety (Rajpurkar et al., 2024), including sharpness, contrast adequacy, noise character-
 058 istics, artifact severity, etc. When these quality attributes are compromised, the resulting subopti-
 059 mal images can lead to diagnostic errors, missed pathologies, or erroneous clinical interpretations,
 060 potentially causing severe patient harm and undermining the integrity of clinical decision-making
 061 processes (Blackmore et al., 2011).

062 Current medical IQA approaches predomi-
 063 nantly produce scalar scores using (1) no-
 064 reference methods (Xun et al., 2025; Herath
 065 et al., 2025), which infer perceptual quality of
 066 an image through statistical feature extraction
 067 without a reference, and (2) full-reference simi-
 068 larity metrics such as PSNR, SSIM (Hore &
 069 Ziou, 2010), and LPIPS (Zhang et al., 2018).
 070 These methods provide standardized evaluation
 071 metrics and enable automated IQA. However,
 072 they exhibit the following fundamental defi-
 073 ciencies. (1) *Poor generalization* (Herath et al.,
 074 2025). Medical image quality is influenced by
 075 complex and heterogeneous factors, including
 076 noise characteristics, contrast adequacy, arti-
 077 fact severity, distortion patterns, and modality-
 078 specific degradations across diverse imaging
 079 modalities such as magnetic resonance imag-
 080 ing (MRI), computed tomography (CT), fundus
 081 photography, histopathology, and endoscopy.
 082 Yet, existing methods typically rely on simple
 083 regression models (Su et al., 2023) or hand-
 084 crafted statistical indices (Dohmen et al., 2024),
 085 which are ill-suited to capture this breadth of quality-affecting factors. As a result, they tend to gen-
 086 eralize poorly to unseen distortions, new modalities or scanners, and different imaging protocols.
 087 (2) *Lack of human-like reasoning process for result interpretation*. Most methods produce scalar
 088 IQA scores, which do not fully reflect the causes of image quality degradation and may be un-
 089 reasonable in certain cases. For instance, as illustrated in Figure 2, when evaluating two medical
 090 images, clinicians typically identify specific degradations first (e.g., metal streak artifacts in Image
 091 #1 vs. reconstruction blurring in Image #2) before assessing their clinical impact. Despite metal arti-
 092 fifacts, Image #1 preserves clear anatomical boundaries and sharp tissue visualization, while Image
 093 #2 suffers from texture loss and unnatural intensity variations. Consequently, Image #1 provides
 094 better clinical quality. However, traditional score-based metrics often favor the smoother Image
 095 #2, contradicting human judgment. Such perceptual reasoning requires understanding the clinical
 096 significance of quality factors, which current automated approaches cannot effectively capture.

097 Recent advances in multimodal large language models (MLLMs) have shown promising capabilities
 098 in medical visual reasoning tasks (OpenAI, 2023; Liu et al., 2024; Dai et al., 2024; Saab et al.,
 099 2024; Su et al., 2025). Theoretically, MLLMs could potentially address existing IQA challenges by
 100 decomposing quality assessment into interconnected subtasks: degradation identification, severity
 101 quantification, clinical impact analysis, and comparative reasoning. Unlike traditional approaches
 102 that yield opaque scores, MLLM-based assessment can provide explicit chains of thought (Wu et al.,
 103 2024a; You et al., 2024), offering interpretable and clinically meaningful evaluations. However,
 104 critical questions remain unanswered about MLLMs’ actual capabilities in medical IQA: Can they
 105 truly generalize to the fine-grained, diverse, and complex quality factors spanning different imaging
 106 modalities? Do they possess genuine reasoning abilities to understand the clinical significance of
 107 various degradations? Existing MLLM evaluation frameworks focus mainly on natural images (Wu
 108 et al., 2024b) or high-level medical semantics (Ye et al., 2024), lacking systematic benchmarks
 109 that assess quality-related perceptual and reasoning skills across diverse medical modalities. This
 110 absence of specialized benchmarks has been a major barrier to developing and validating effective
 111 frameworks.

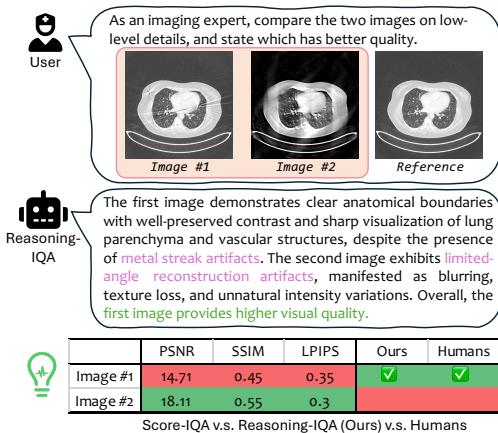


Figure 2: Comparison of Reasoning IQA with score-based IQA. Unlike purely numerical scores, Reasoning IQA identifies distortion types and their relative impact, yielding results more consistent with human judgment.

No-reference Scenario		
Perception		
Yes-or-no	<p>Q: Is the optic disc fully visible and clearly delineated in the image? O: A:Yes B:<input checked="" type="checkbox"/> No</p>	<p>Q: Is the anterior posterior axis of the aorta clearly visible in the MRI image? O: A:Yes B:<input checked="" type="checkbox"/> No</p>
What	<p>Q: What is the primary artifact visible in the CT image? O: A:Incomplete projections <input checked="" type="checkbox"/> B:Motion blur C:Grid artifacts</p>	<p>Q: What type of artifacts are predominantly visible in this MRI scan? O: A: Motion artifacts B:<input checked="" type="checkbox"/> Undersampling artifacts C: Metal artifacts D:Noise</p>
How	<p>Q: How would you rate the overall quality of this image? O: A:Good B:Usable C:Reject <input checked="" type="checkbox"/></p>	<p>Q: How would you rate the visibility of vessel branching detail in the image? O: A:Good B:<input checked="" type="checkbox"/> Usable C:Reject</p>
Reasoning		
	<p>Authentic artifact</p>	<p>The image shows the gastrointestinal tract with visible mucosal folds and a central lumen. Globally, the image has moderate clarity and adequate contrast but suffers from uneven illumination. Locally, specular reflections and artifacts are present in both central and peripheral regions. Some areas also exhibit blur and saturation. Overall, the quality of this image is usable.</p>
	<p>AI-generated</p>	<p>This fundus image shows a wide field of view with the optic disc, macula, and retinal vessels under a yellow-green tone. The clarity and contrast are acceptable, with relatively uniform illumination, though the yellow cast slightly affects detail perception. Key structures like the optic disc and major vessels are visible. Overall, the quality of this image is good.</p>
	<p>Simulate Artifact</p>	<p>This is an MRI image of the knee. The image exhibits low clarity and contrast, with a high noise level and poor overall resolution due to significant undersampling. Undersampling artifacts, such as aliasing and streaking, are prominent, severely impacting the ability to accurately interpret anatomical structures. Overall, the quality of this image is reject.</p>
Paired Scenario		
	<p>Coarse-grained</p> <p>The first image shows excellent clarity and sharpness with minimal motion artifacts, enabling clear visualization of anatomical structures. By contrast, the second image has severe motion artifacts, including blurring and ghosting, which obscure details and reduce diagnostic quality. Overall, the first image is of higher quality.</p>	<p>Fine-grained</p> <p>The first image demonstrates clear anatomical boundaries with well-preserved contrast and sharp visualization of lung parenchyma and vascular structures, despite the presence of metal streak artifacts. The second image exhibits sparse view reconstruction artifacts, manifested as blurring, texture loss, and unnatural intensity variations. Overall, the first image provides higher visual quality.</p>

Figure 3: Examples of question types in MedQ-Bench, covering MCQA perception tasks (Yes-No / What / How), open-ended reasoning, and pair/multi-image comparison.

To bridge the gap between existing medical IQA methods, we propose a novel *perception–reasoning paradigm*. This paradigm mirrors clinicians’ cognitive workflow: first perceiving quality-related attributes in images, assessing their severity, and evaluating their potential impact on clinical diagnosis, and then making overall quality judgments through logical reasoning. Building on this paradigm, we introduce **MedQ-Bench**, the first comprehensive benchmark that systematically evaluates the medical IQA capabilities of MLLMs. Our primary contributions are as follows:

- **Pioneering evaluation framework for medical image quality assessment.** **MedQ-Bench** introduces a systematic evaluation methodology that comprehensively assesses both quality-based perceptual and reasoning capabilities for MLLMs. The framework extends beyond traditional IQA scoring to incorporate quality-related perception assessment, fine-grained comparative analysis, and quality-aware reasoning evaluation. The protocol supports both no-reference and full-reference paradigms, enabling systematic assessment ranging from coarse-grained to fine-grained perceptual discrimination tasks.
- **Multi-dimensional judging protocol with human–AI alignment validation.** To evaluate reasoning ability, we design a multi-dimensional judging protocol that scores model outputs along four complementary axes. We further perform rigorous human–AI alignment validation by comparing our LLM-based evaluations with radiologists, demonstrating the reliability of the proposed evaluation framework.
- **Comprehensive, clinically representative, multi-source dataset.** Covering 5 imaging modalities and 40+ quality attributes, **MedQ-Bench** blends authentic clinical images, simulated degraded images via physics-based reconstruction, and AI-generated images to encompass diverse real-world and controlled scenarios. This comprehensive dataset enables robust evaluation across both realistic clinical conditions and challenging scenarios.
- **Comprehensive empirical analysis.** We conduct extensive evaluations of state-of-the-art MLLMs, spanning open-source and commercial systems, both general-purpose and medical-specialized. Our systematic analysis reveals significant performance gaps in modality-specific perception capabilities, underscoring the need for targeted improvements for clinical readiness.

162 2 CONSTRUCTING THE MEDQ-BENCH
163

164 2.1 BENCHMARK SCOPE AND MODALITIES
165

166 Clinical image quality is fundamental to diagnostic reliability, yet existing evaluation methods rely
167 primarily on score-based metrics that overlook the comprehensive assessment of image quality per-
168 ception and reasoning capabilities. MedQ-Bench is specifically designed to systematically eval-
169 uate the visual quality perception and reasoning capabilities of multimodal large language models
170 (MLLMs) within the medical imaging domain. Let $\mathcal{M} = \{M_1, M_2, \dots, M_5\}$ represent the set of
171 five medical imaging modalities, where each modality M_i is associated with a distinct set of quality
172 attributes $\mathcal{A}_i = \{a_{i,1}, a_{i,2}, \dots, a_{i,k}\}$. The quality assessment task can be formulated as learning a
173 mapping function $f : \mathcal{I} \times \mathcal{Q} \rightarrow \mathcal{R}$ that takes an image $I \in \mathcal{I}$ and question $q \in \mathcal{Q}$ as input and
174 produces a response $r \in \mathcal{R}$.

175 To capture the diversity and complexity of real-world clinical imaging, MedQ-Bench encompasses
176 five representative modalities: Magnetic Resonance Imaging (MRI), Computed Tomography (CT),
177 endoscopy, histopathology imaging, and fundus photography. Let $\mathcal{D}_i = \{d_{i,1}, d_{i,2}, \dots, d_{i,n}\}$ denote
178 the set of degradation types specific to modality M_i . Each modality exhibits distinct degradation
179 characteristics due to its physical acquisition principles, where degradations can be modeled as
180 transformations $T_d : \mathcal{I} \rightarrow \mathcal{I}'$ that modify the original image I based on degradation type $d \in \mathcal{D}_i$.
181 For instance, MRI is particularly susceptible to motion and magnetic susceptibility artifacts, and CT
182 is prone to low-dose noise and metal-induced streak artifacts. This multi-modality design ensures
183 that the benchmark reflects the broad spectrum of perceptual challenges encountered in practice.

184 For each modality, MedQ-Bench incorporates images from three complementary sources: authentic
185 clinical images containing naturally occurring artifacts; synthetically degraded images that replicate
186 modality-specific distortions in a controlled manner; and AI-generated or reconstructed images pro-
187 duced by enhancement, translation, or reconstruction models, which may introduce hallucinations or
188 subtle structural inconsistencies. Let $\mathcal{S} = \{\mathcal{S}_{\text{real}}, \mathcal{S}_{\text{synth}}, \mathcal{S}_{\text{AI}}\}$ denote the three image sources, where
189 each source \mathcal{S}_k contributes a subset of images with specific degradation characteristics $\mathcal{D}_k \subseteq \bigcup_i \mathcal{D}_i$.
190 This tri-source strategy enables the benchmark to cover both naturally occurring degradations and
191 algorithm-induced artifacts, ensuring a balanced evaluation of MLLM robustness across real-world
192 and algorithmic distortion scenarios.

193 2.2 BENCHMARK ON IQA PERCEPTION ABILITY
194

195 Before evaluating sophisticated reasoning capabilities, it is essential to establish whether MLLMs
196 possess fundamental perceptual abilities to recognize basic image quality attributes.
197

198 2.2.1 QUESTION TYPES
199

200 The perception-focused MCQA setting evaluates direct visual perception using single-image
201 prompts, without requiring domain-specific diagnostic reasoning. These tasks represent the most
202 basic level of quality assessment capability, asking models to simply identify “what they see” rather
203 than explain “why they see it.” For each image, three canonical subtypes of questions are in-
204 cluded: (1) **Yes-or-No**: Binary classification tasks where $\mathcal{R}_{\text{YN}} = \{0, 1\}$ and the model predicts
205 $\hat{y} = \arg \max_{y \in \{0,1\}} P(y | I, q)$. Examples include “Is this image clear?” or “Does this image
206 contain artifacts?” (2) **What**: Multi-class identification tasks where $\mathcal{R}_{\text{What}} = \{c_1, c_2, \dots, c_K\}$ rep-
207 resents K possible degradation types, and the model selects $\hat{c} = \arg \max_{c \in \mathcal{R}_{\text{What}}} P(c | I, q)$. These
208 tasks ask models to identify specific types of artifacts or degradations present in the image. (3)
209 **How**: Severity assessment tasks where $\mathcal{R}_{\text{How}} = \{s_1, s_2, \dots, s_L\}$ represents L severity levels, and
210 the model predicts $\hat{s} = \arg \max_{s \in \mathcal{R}_{\text{How}}} P(s | I, q)$. These tasks evaluate the model’s ability to assess
211 the degree or intensity of observed quality issues.

212 2.2.2 QUADRANTS FOR LOW-LEVEL VISUAL CONCERNs
213

214 **Axis 1: No Degradation vs Degradation Severity Levels.** The primary axis differentiates medical
215 images based on their quality degradation status: 1) **No Degradation** refers to medical images that
 maintain optimal quality standards without artifacts or distortions, and 2) degradation with Severity

216 Levels encompasses images with varying degrees of quality issues, further subdivided into *mild*
217 *Degradation* and *severe Degradation*.

218 **Axis 2: General Medical Questions vs Modality-specific Questions.** Quality perception in medical
219 imaging intertwines with modality-specific technical characteristics. For instance, motion artifacts
220 manifest differently in MRI versus CT scans. We curate *modality-specific questions* that
221 require understanding unique technical characteristics of specific imaging modalities (e.g., “Does
222 this MRI show susceptibility artifacts?”), while *general medical questions* focus on universal quality
223 concepts applicable across modalities (e.g., “Is this image clear?”). This distinction evaluates
224 both fundamental quality perception and specialized modality knowledge.

225

226 2.3 BENCHMARK ON IQA REASONING ABILITY

227

228 2.3.1 NO-REFERENCE REASONING TASKS.

229 While MCQA constrains answers to predefined choices, reasoning tasks assess a model’s ability to
230 autonomously describe and explain quality-related observations in natural language. These tasks
231 require generating comprehensive responses $w_{1:T} = \{w_1, w_2, \dots, w_T\}$ that systematically detail
232 multiple aspects of image quality assessment: (1) modality and anatomical region identification; (2)
233 specific quality degradation characterization including type and severity; (3) technical attribution of
234 underlying causes; (4) assessment of diagnostic impact and clinical implications; and (5) definitive
235 quality judgment with good/usable/reject recommendation. The reasoning tasks evaluate whether
236 models can perform structured quality analysis that mirrors expert clinical assessment, moving be-
237 yond simple classification to demonstrate understanding of the relationship between technical image
238 properties, degradation mechanisms, and clinical utility.

239

240 2.3.2 COMPARISON REASONING TASKS.

241 Many clinical workflows require comparative quality assessment between two versions of the same
242 study, such as “original vs. reconstructed” or outputs from competing reconstruction algorithms.
243 For image pairs (I_A, I_B) , the comparative task seeks to determine preference $P(I_A \succ I_B)$ based
244 on overall quality assessment. Models must identify which image exhibits higher diagnostic quality
245 and provide detailed explanations for their judgment, such as explaining why one reconstruction
246 algorithm preserves anatomical detail better than another.

247 Comparative tasks are further categorized by the perceptual gap between images. 1) *Coarse-grained*
248 comparisons involve clearly visible quality differences, making them relatively straightforward for
249 both humans and models. 2) *Fine-grained* comparisons involve subtle differences in noise patterns,
250 contrast, or structure fidelity, requiring heightened sensitivity to nuanced quality cues that may only
251 be apparent upon careful inspection. This design enables separate evaluation of basic discrimination
252 ability and advanced perceptual subtlety that approaches expert-level assessment sensitivity.

253

254 2.3.3 EVALUATION METRICS

255 **Multi-dimensional judging protocol** The reasoning tasks require more nuanced evaluation ap-
256 proaches due to their subjective nature and the complexity of natural language responses. Re-
257 cent studies have demonstrated GPT-4o to be a reliable evaluation tool for complex reasoning
258 tasks. We assess model outputs \mathcal{O} across four complementary dimensions, each scored on a
259 discrete scale $s \in \{0, 1, 2\}$: **(1) Completeness.** $C(\mathcal{O}, \mathcal{R}) = \frac{1}{|\mathcal{K}_{\mathcal{R}}|} \sum_{k \in \mathcal{K}_{\mathcal{R}}} \mathbb{I}[k \in \mathcal{K}_{\mathcal{O}}]$ mea-
260 sures the coverage of key visual information from the reference description \mathcal{R} , where $\mathcal{K}_{\mathcal{R}}$ and
261 $\mathcal{K}_{\mathcal{O}}$ represent the sets of key visual information in reference and output respectively. Higher
262 scores indicate more comprehensive description of observable quality issues. **(2) Precise-
263 ness.** $P(\mathcal{O}, \mathcal{R}) = 1 - \frac{1}{|\mathcal{K}_{\mathcal{O}}|} \sum_{k \in \mathcal{K}_{\mathcal{O}}} \mathbb{I}[\text{contradict}(k, \mathcal{R})]$ quantifies consistency between model
264 output and reference by penalizing semantic contradictions. **(3) Consistency.** $S(\mathcal{O}, \mathcal{R}) =$
265 $f_{\text{consistency}}(\text{reasoning}(\mathcal{O}), \text{conclusion}(\mathcal{O}), \mathcal{R})$ evaluates the internal logical consistency between the
266 reasoning path $\text{reasoning}(\mathcal{O})$ and the final quality judgment $\text{conclusion}(\mathcal{O})$, where $f_{\text{consistency}}$ re-
267 turns a score based on logical coherence assessment. **(4) Quality Accuracy.** $Q(\mathcal{O}, \mathcal{R}) =$
268 $\mathbb{I}[\text{comparison}(\mathcal{O}) = \text{comparison}(\mathcal{R})]$ assesses whether the final quality comparison judgment cor-
269 rectly identifies which image has higher quality, matching the reference assessment. This binary
metric focuses on the correctness of the ultimate quality decision.

270 **Human–AI Alignment Validation** To ensure the reliability and validity of our automated evalua-
 271 tion, we conducted a rigorous alignment validation between GPT-4o judgments and expert assess-
 272 ments. A total of 200 cases were randomly sampled from the development dataset and independently
 273 evaluated by three board-certified medical imaging specialists under a double-blinded protocol.

274 For human–AI alignment, we employed quadratic weighted Cohen’s kappa (Cohen, 1968) for ordi-
 275 nal ratings:

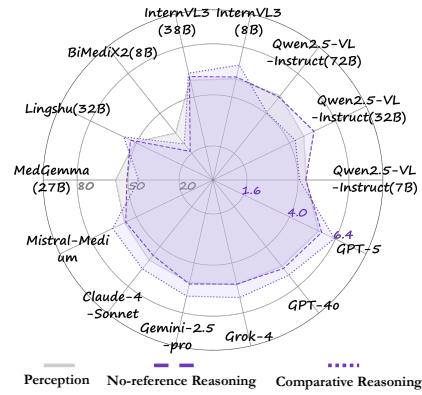
$$\kappa_w = 1 - \frac{\sum_{i,j} w_{ij} O_{ij}}{\sum_{i,j} w_{ij} E_{ij}}, \quad (1)$$

279 where O_{ij} is the observed agreement matrix, E_{ij} the expected agreement matrix, and $w_{ij} = \frac{(i-j)^2}{(k-1)^2}$
 280 the quadratic weights penalizing larger disagreements more severely. We further conducted iter-
 281 ative prompt refinement to maximize concordance between GPT-4o and expert consensus. Final
 282 alignment results are reported in Section 3.4.
 283

284 3 RESULTS

287 To investigate MLLMs’ image quality per-
 288 ception ability, we present a comprehensive
 289 evaluation of MedQ-Bench across 14 up-to-
 290 date popular MLLMs under zero-shot settings.
 291 We evaluate these 14 multimodal large lan-
 292 guage models across three categories: open-
 293 source MLLMs (Qwen2.5-VL-Instruct vari-
 294 ants (Wang et al., 2024a), InternVL3 mod-
 295 els (Chen et al., 2024b)), medical-specialized
 296 MLLMs (BiMediX2 (Peng et al., 2024), Ling-
 297 shu (Wang et al., 2024b), MedGemma (Saab
 298 et al., 2024)), and commercial systems (GPT-
 299 5 (OpenAI, 2024b), GPT-4o (OpenAI, 2024a),
 300 Gemini-2.5-Pro (Reid et al., 2024), Grok-
 301 4 (xAI Team, 2024), Claude-4-Sonnet (An-
 302 thropic, 2024), Mistral-Medium-3 (Jiang et al.,
 303 2023)).

3.1 FINDINGS ON PERCEPTION



304 Figure 4: Overall Performance Results

Sub-categories	Perception				Reasoning					
	Model (variant)	Yes-or-No↑	What↑	How↑	Overall↑	Comp.↑	Prec.↑	Cons.↑	Qual.↑	Overall↑
random guess	50.00%	28.48%	33.30%	37.94%						
Non-experts	67.50%	57.50%	57.50%	62.50%	-	-	-	-	-	
Human experts	88.50%	77.50%	77.50%	82.50%	-	-	-	-	-	
Qwen2.5-VL-Instruct (7B)	57.89%	48.45%	54.40%	54.71%	0.715	0.670	1.855	1.127	4.367	
Qwen2.5-VL-Instruct (32B)	67.38%	43.02%	58.69%	59.31%	1.077	0.928	1.977	1.290	<u>5.272</u>	
InternVL3 (8B)	72.04%	47.67%	52.97%	60.08%	0.928	0.878	1.858	1.317	4.983	
InternVL3 (38B)	69.71%	57.36%	52.97%	61.00%	0.964	0.824	<u>1.860</u>	1.317	4.965	
Qwen2.5-VL-Instruct (72B)	78.67%	42.25%	56.44%	63.14%	0.905	0.860	1.896	1.321	4.982	
BiMediX2 (8B)	44.98%	27.52%	27.81%	35.10%	0.376	0.394	0.281	0.670	1.721	
Lingshu (32B)	50.36%	50.39%	51.74%	50.88%	0.624	0.697	1.932	1.059	4.312	
MedGemma (27B)	67.03%	48.06%	50.72%	57.16%	0.742	0.471	1.579	1.262	4.054	
Mistral-Medium-3	65.95%	48.84%	52.97%	57.70%	0.923	0.729	1.566	1.339	4.557	
Claude-4-Sonnet	71.51%	46.51%	54.60%	60.23%	0.742	0.633	1.778	1.376	4.529	
Gemini-2.5-Pro	75.13%	<u>55.02%</u>	50.54%	61.88%	0.878	<u>0.891</u>	1.688	1.561	5.018	
Grok-4	73.30%	48.84%	59.10%	63.14%	<u>0.982</u>	0.846	1.801	1.389	5.017	
GPT-4o	78.48%	49.64%	57.32%	64.79%	<u>1.009</u>	1.027	1.878	1.407	5.321	
GPT-5	82.26%	60.47%	<u>58.28%</u>	68.97%	1.195	1.118	1.837	<u>1.529</u>	5.679	

321 Table 1: Performance of different models on the MCQA perception and reasoning tasks. First place
 322 in each column is bolded; second and third places are underlined. Random guess / Non-experts /
 323 Human experts are excluded from ranking.

To ensure rigorous and unbiased evaluation, the **MedQ-Perception** is equally divided into `dev` (Table 6, for prompt refinement) and `test` (Table 1, for final evaluation) subsets.

Conclusion 1. Clear performance hierarchy emerges across model categories: Our analysis reveals that most MLLMs perform above random guessing across all sub-tasks, indicating promising potential for domain generalization. The results demonstrate a clear performance hierarchy: closed-source frontier models achieve the highest scores, with GPT-5 leading at 68.97% on the test set. Among open-source models, Qwen2.5-VL-Instruct (72B) achieves the best performance at 63.14%, outperforming most commercial models, while *the best medical-specialized models underperform expectations*, with MedGemma (27B) achieving only 57.16%. More details are in Section A.5.1.

Insufficiency 1. Substantial human-AI performance gap remains: Another key finding emerges from our comparison with human performance, where we include both **human experts** (medical imaging technicians and medical imaging PhDs) and **non-experts** as reference points. The best AI model (GPT-5) significantly underperforms human experts (68.97% vs. 82.50%, a gap of 13.53%), yet outperforms non-experts by 6.47%. Given that these models have not undergone specialized training for medical image quality assessment, this suggests substantial potential for improvement in these MLLMs through further fine-tuning.

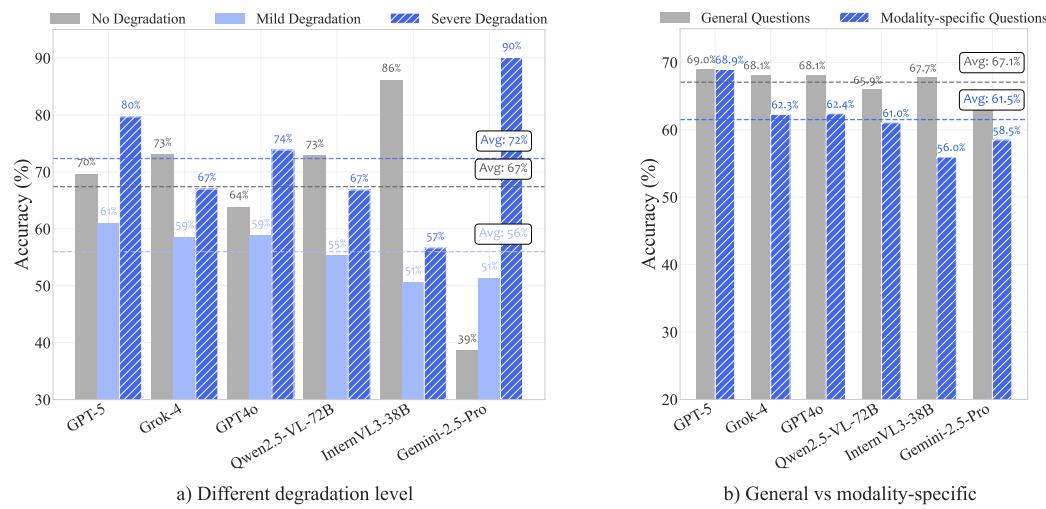


Figure 5: Performance analysis of MLLMs across different evaluation dimensions. (a) Different degradation level performance . (b) General vs modality-specific question.

Insufficiency 2. The LVLMs are not robust among different perceptual types: Task-specific analysis reveals distinct patterns across different evaluation dimensions. Performance analysis across different degradation levels (Figure 5(a)) demonstrates that mild degradation represents the most challenging detection scenario, with average accuracy dropping to 56% compared to 72% for no degradation and 67% for severe degradation. This indicates that subtle quality issues are harder to identify than obvious artifacts. Top-performing models like GPT-5 demonstrate a degree of consistency in performance across degradation levels. We further investigate the difference between general and modality-specific medical questions. As shown in Figure 5(b), most models perform better on general questions than on modality-specific tasks, whereas GPT-5 demonstrates the most balanced performance across question types. This suggests that robust medical image quality assessment requires specialized understanding of modality-specific visual features.

3.2 FINDINGS ON NO-REFERENCE REASONING

Conclusion 2. Limited low-level visual reasoning capabilities across all models: For no-reference reasoning capabilities (Table 1), GPT-5 still demonstrates the best performance, particularly excelling in the relevance dimension. However, even the most advanced MLLMs fail to achieve excellent scores in completeness and preciseness, with the highest scores being only 1.293/2.0 for

Model	Comp. \uparrow	Prec. \uparrow	Cons. \uparrow	Qual. \uparrow	Overall \uparrow
Qwen2.5-VL-7B	0.714	0.902	1.316	1.143	4.075
Qwen2.5-VL-32B	0.692	0.752	1.895	0.962	4.301
Qwen2.5-VL-72B	0.737	0.977	1.233	1.113	4.060
InternVL3-8B	0.985	1.278	1.797	1.474	5.534
InternVL3-38B	1.075	1.083	1.571	1.414	5.143
BiMediX2-8B	0.474	0.549	0.639	0.511	2.173
MedGemma-27B	0.684	0.692	1.128	1.000	3.504
Lingshu-32B	0.729	1.015	1.586	1.323	4.653
Mistral-Medium-3	0.872	1.203	1.827	1.338	5.240
Claude-4-Sonnet	0.857	1.083	1.910	1.481	5.331
Gemini-2.5-Pro	1.053	1.233	1.774	1.534	5.594
Grok-4	1.150	1.233	1.820	1.459	5.662
GPT-4o	1.105	1.414	1.632	1.562	5.713
GPT-5	1.293	1.556	1.925	1.564	6.338

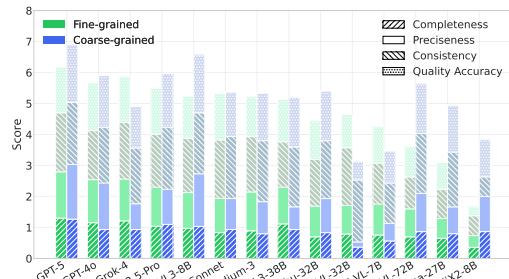


Figure 6: Comparative reasoning performance analysis. Left: Detailed performance scores across three evaluation dimensions for all models. Right: Visual comparison of overall performance patterns across model categories.

completeness and 1.556/2.0 for preciseness. In general, most models only reach an acceptable baseline level. Current MLLM models possess relatively limited and elementary low-level visual reasoning abilities, struggling to provide complete and accurate descriptions of low-level visual information. The consistently high consistency scores indicate that most MLLMs can follow abstract instructions reasonably well, suggesting that the main bottleneck for improving MLLM descriptive capabilities lies in the perception of low-level attributes rather than instruction following.

3.3 FINDINGS ON COMPARISON REASONING

Insufficiency 3. Paired comparison reveals fundamental limitations in fine-grained analysis: Paired image comparison tasks pose the greatest challenge to current multimodal large language models (MLLMs), requiring models to perform fine-grained quality comparisons between similar images that may only differ by varying degrees. We evaluate model performance across two difficulty levels: fine-grained differences and coarse-grained differences. Figure 6 (right) presents detailed performance analysis across different difficulty levels, with more complete tabular results available in Table 9 in the appendix. Overall, most models perform better under coarse-grained differences, while a few models, such as Grok-4 and Qwen2.5-VL-7B/32B, perform better under fine-grained differences but lose performance on coarse-grained tasks. Among them, GPT-5 achieved the highest overall score, while medical-specialized models such as BiMediX2 showed notably insufficient performance.

3.4 HUMAN-AI ALIGNMENT VALIDATION

Strong human-AI alignment validates our evaluation framework: To validate the reliability of our automated evaluation approach, we conducted a comprehensive human-AI alignment study comparing human expert assessments with GPT-4o automated scoring. We evaluated 200 randomly sampled image quality assessments across three key dimensions: completeness, preciseness, and consistency. The confusion matrices in the appendix (Figure 13) demonstrate strong alignment between human expert scores and GPT-4o automated evaluation across all three dimensions, with consistently high accuracy rates: 83.3% for completeness, 87.0% for preciseness, and 90.5% for consistency, with all individual class recall rates exceeding 80%.

These results validate that our automated quality assessment system achieves strong alignment with human expert judgment across all evaluation dimensions, with high accuracy rates demonstrating that our evaluation framework can serve as a reliable substitute for human evaluation. Beyond accuracy, we further assessed inter-rater agreement using quadratic weighted Cohen’s κ_w (Table 10), achieving consistently high values (0.774–0.985) that confirm substantial agreement beyond chance and validate our framework as a reliable surrogate for large-scale human evaluation.

432 4 RELATED WORK 433

434 **Medical Multimodal Large Language Models and Benchmarks.** Multimodal Large Language
435 Models (MLLMs) have demonstrated remarkable capabilities in understanding and reasoning about
436 visual content through natural language. General-purpose models like GPT-4V (OpenAI, 2023),
437 LLaVA (Liu et al., 2024), and Qwen-VL (Wang et al., 2024a) have shown strong performance across
438 diverse vision-language tasks. To address healthcare-specific requirements, medical-specialized
439 variants such as (Wang et al., 2024b; Saab et al., 2024; Peng et al., 2024; Su et al., 2025; Xu et al.,
440 2025) have emerged through domain-targeted pretraining and alignment. Recent medical bench-
441 marks have been developed to evaluate these models systematically, including (Ye et al., 2024),
442 which provides comprehensive multimodal evaluation for general medical AI. However, existing
443 medical benchmarks focus primarily on high-level diagnostic tasks rather than low-level perceptual
444 quality assessment (Chen et al., 2024a).

445 **Score-based Image Quality Assessment.** Traditional image quality assessment methods produce
446 numerical scores to quantify image quality, categorized into No-Reference (NR), Full-Reference
447 (FR), and Reduced-Reference approaches. NR methods like BRISQUE (Mittal et al., 2012),
448 NIQE (Zhang et al., 2015), and deep learning approaches including CNNIQA (Kang et al., 2014)
449 and MUSIQ (Ke et al., 2021) assess quality without reference images. FR methods compare against
450 pristine references using metrics like PSNR, SSIM (Wang et al., 2004), VIF (Sheikh & Bovik,
451 2006), and learned perceptual metrics like LPIPS (Zhang et al., 2018). Recent advances include
452 transformer-based approaches like TReS (Golestaneh et al., 2022) and quality-aware pretraining
453 methods. However, these methods yield only scalar scores, offering limited interpretability regard-
454 ing specific quality factors, and such technical measures often show weak alignment with clinical
455 workflows (Zhang et al., 2024; Blackmore et al., 2011).

456 **MLLM-based Image Quality Assessment.** Recent advances have introduced multimodal lan-
457 guage models for image quality assessment (IQA), which enable more interpretable and reasoning-
458 based evaluation. For example, Q-Instruct (Wu et al., 2024a) and DepictQA (You et al., 2024)
459 generate natural language descriptions of quality factors, while Q-Bench (Wu et al., 2024b) offers a
460 systematic framework for evaluating low-level vision tasks. Building on this line, IQAGPT (Chen
461 et al., 2023) integrates vision-language models with ChatGPT for CT image quality assessment,
462 showing the feasibility of producing both quality scores and textual reports. However, its scope is
463 limited to CT images and remains focused on score prediction rather than comprehensive reasoning.
464 Likewise, Ultrasound-QBench (Miao et al., 2025) provides evaluation for ultrasound imaging but
465 restricts tasks to classification and scoring within a single modality.

466 5 CONCLUSION 467

468 We introduced **MedQ-Bench**, the first benchmark to systematically evaluate medical image quality
469 assessment (IQA) capabilities of multimodal large language models through a perception–reasoning
470 paradigm. Unlike conventional score-based metrics, **MedQ-Bench** jointly assesses quality-related
471 perception and reasoning across five imaging modalities and more than forty degradation types via
472 three complementary tracks: perception tasks, no-reference reasoning, and paired comparison rea-
473 soning. Our large-scale zero-shot evaluation of 14 state-of-the-art MLLMs, including open-source,
474 medical-specialized, and commercial systems, yields several key findings. Substantial performance
475 gaps remain between AI models and human experts, particularly in detecting subtle degradations
476 critical to clinical practice. Current models exhibit preliminary but unstable perceptual and reason-
477 ing abilities, often failing to produce complete and precise quality descriptions. Medical-specialized
478 models unexpectedly underperform general-purpose ones, calling into question the effectiveness of
479 current domain adaptation strategies. Moreover, models show marked weaknesses in fine-grained
480 comparisons and mild degradation detection, precisely where reliable quality control is most needed.
481 By moving beyond high-level diagnostic reasoning toward foundational quality perceptual and rea-
482 soning skills, **MedQ-Bench** establishes a clinically grounded and interpretable standard for mea-
483 suring and advancing medical IQA. We anticipate that it will inform the development of MLLMs
484 with stronger low-level visual understanding and trustworthy reasoning, paving the way for safe and
485 reliable integration of automated quality control into clinical imaging workflows.

486 REFERENCES

487

- 488 AAPM CT-MAR Challenge Organizers. The aapm ct metal artifact reduction (ct-mar)
489 grand challenge — dataset. <https://www.aapm.org/GrandChallenge/CT-MAR/>; <https://qtim-challenges.southcentralus.cloudapp.azure.com/competitions/1>, 2023. URL <https://www.aapm.org/GrandChallenge/CT-MAR/>. Generated using XCIST, with hybrid data sim-
490 ulation framework combining clinical images and virtual metal objects.
491
492
- 493 Anthropic. Claude 4: Constitutional ai with harmlessness from ai feedback. *Technical Report*, 2024.
494 Available at <https://www.anthropic.com/clause>.
495
496 Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse
497 problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
498
499 Craig C Blackmore, Robert S Mecklenburg, and Gary S Kaplan. Evidence-based radiology: a new
500 approach to the practice of radiology. *Radiology*, 259(3):615–620, 2011.
501
502 Tianhe Chen, Shuai Liu, Yizhou Zhang, and Rui Zhao. A comprehensive study of multimodal large
503 language models for image quality assessment. pp. 143–159, 2024a.
504
505 Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong
506 Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning
507 for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer*
508 *vision and pattern recognition*, pp. 24185–24198, 2024b.
509
510 Zhihao Chen, Bin Hu, Chuang Niu, Tao Chen, Yuxin Li, Hongming Shan, and Ge Wang.
511 Iqagpt: Image quality assessment with vision-language and chatgpt models. *arXiv preprint*
512 *arXiv:2312.15663*, 2023.
513
514 Jacob Cohen. Weighted kappa: Nominal scale agreement provision for scaled disagreement or
515 partial credit. *Psychological bulletin*, 70(4):213, 1968.
516
517 Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang,
518 Boyang Li, Pascal N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-
519 language models with instruction tuning. *Advances in Neural Information Processing Systems*,
520 36:49250–49267, 2024.
521
522 Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao
523 Yu, Xiaonan Nie, Ziang Song, et al. Emerging properties in unified multimodal pretraining. *arXiv*
524 *preprint arXiv:2505.14683*, 2025.
525
526 Melanie Dohmen, Tuan Truong, Ivo M Baltruschat, and Matthias Lenga. Five pitfalls when assess-
527 ing synthetic medical images with reference metrics. In *MICCAI Workshop on Deep Generative*
528 *Models*, pp. 150–159. Springer, 2024.
529
530 Huazhu Fu, Boyang Wang, Jianbing Shen, Shanshan Cui, Yanwu Xu, Jiang Liu, and Ling Shao.
531 Evaluation of retinal image quality assessment networks in different color-spaces. In *International*
532 *conference on medical image computing and computer-assisted intervention*, pp. 48–56. Springer,
533 2019.
534
535 Moritz Fuchs, Ssharvien Kumar R Sivakumar, Mirko Schöber, Niklas Woltering, Marie-Lisa Eich,
536 Leonille Schweizer, and Anirban Mukhopadhyay. Harp: Unsupervised histopathology artifact
537 restoration. In *Medical Imaging with Deep Learning*, 2024.
538
539 S. Alireza Golestaneh, Saba Dadsetan, and Kris M. Kitani. No-reference image quality assessment
540 via transformers, relative ranking, and self-consistency. In *Proceedings of the IEEE/CVF Winter*
541 *Conference on Applications of Computer Vision*, pp. 1220–1230, 2022.
542
543 Matthieu Guerquin-Kern, M Haberlin, Klaas Paul Pruessmann, and Michael Unser. A fast wavelet-
544 based reconstruction method for magnetic resonance imaging. *IEEE transactions on medical*
545 *imaging*, 30(9):1649–1660, 2011.
546
547 HMSS Herath, HMKKMB Herath, Nuwan Madusanka, and Byeong-II Lee. A systematic review of
548 medical image quality assessment. *Journal of Imaging*, 11(4):100, 2025.

-
- 540 Alain Hore and Djamel Ziou. Image quality metrics: Psnr vs. ssim. In *2010 20th international*
541 *conference on pattern recognition*, pp. 2366–2369. IEEE, 2010.
- 542
- 543 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chap-
544 lot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,
545 Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril,
546 Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL <https://arxiv.org/abs/2310.06825>.
- 547
- 548 Avinash C Kak and Malcolm Slaney. *Principles of computerized tomographic imaging*. SIAM,
549 2001.
- 550
- 551 Le Kang, Peng Ye, Yi Li, and David Doermann. Convolutional neural networks for no-reference
552 image quality assessment. In *Proceedings of the IEEE conference on computer vision and pattern*
553 *recognition*, pp. 1733–1740, 2014.
- 554
- 555 Neel Kanwal. Histoartifacts, March 2024. URL <https://doi.org/10.5281/zenodo.10809442>.
- 556
- 557 Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. Musiq: Multi-scale image
558 quality transformer. In *Proceedings of the IEEE International Conference on Computer Vision*,
559 pp. 5148–5157, 2021.
- 560
- 561 Mathieu Lamard, Jean-Louis Coatrieux, Philippe Dequidt, Marc Le Berre, Christian Roux, and
562 Basel Solaiman. Checklist for artificial intelligence in medical imaging (claim): 2024 update. In
563 *Radiology: Artificial Intelligence*, volume 6, pp. e240300, 2024.
- 564
- 565 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances*
566 *in Neural Information Processing Systems*, 36:36820–36835, 2024.
- 567
- 568 Hongyi Miao, Jun Jia, Yankun Cao, Yingjie Zhou, Yanwei Jiang, Zhi Liu, and Guangtao Zhai.
Ultrasound-qbench: Can llms aid in quality assessment of ultrasound imaging? *arXiv preprint*
569 *arXiv:2501.02751*, 2025.
- 570
- 571 Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. No-reference image quality assess-
572 ment in the spatial domain. *IEEE Transactions on image processing*, 21(12):4695–4708, 2012.
- 573
- 574 OpenAI. Gpt-4v(ision) system card. *Technical Report*, 2023. URL <https://openai.com/research/gpt-4v-system-card>.
- 575
- 576 OpenAI. Gpt-4 technical report, 2024a. URL <https://arxiv.org/abs/2303.08774>.
- 577
- 578 OpenAI. Gpt-5: Artificial general intelligence is near. *Technical Report*, 2024b. URL <https://openai.com/gpt-5>.
- 579
- 580 Sahal Shaji Peng, Vandan Gorade Narayan, Sreya Deshpande, et al. Bimedix2: Bio-medical expert
581 lmm for diverse medical modalities. *arXiv preprint arXiv:2405.20157*, 2024.
- 582
- 583 Gorkem Polat, Deniz Sen, Alperen Inci, and Alptekin Temizel. Endoscopic artefact detection with
584 ensemble of deep neural networks and false positive elimination. In *Proceedings of the 2nd*
585 *International Workshop and Challenge on Computer Vision in Endoscopy, EndoCV@ISBI 2020,*
586 *Iowa City, Iowa, USA, 3rd April 2020*, volume 2595, pp. 8–12. CEUR-WS.org, 2020.
- 587
- 588 Pranav Rajpurkar, Erin Chen, Oishi Banerjee, and Eric J Topol. Ai in diagnostic imaging: Revolu-
589 tionising accuracy and efficiency. *Clinical Radiology*, 79(2):e132–e140, 2024.
- 590
- 591 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-
592 Baptiste Alayrac, et al. Gemini 2.5: Unlocking multimodal understanding across image, video,
593 and audio. *arXiv preprint arXiv:2312.11805*, 2024.
- 594
- 595 Khaled Saab, Yi Tay, et al. Medgemma: A medical large language model specialized for high-stakes
596 applications. *arXiv preprint arXiv:2409.03278*, 2024.
- 597
- 598 Hamid R Sheikh and Alan C Bovik. Image information and visual quality. *IEEE Transactions on*
599 *Image Processing*, 15(2):430–444, 2006.

- 594 Jialin Su, Meifang Li, Yongping Lin, Liu Xiong, Caixing Yuan, Zhimin Zhou, and Kunlong Yan.
595 Deep learning-driven multi-view multi-task image quality assessment method for chest ct image.
596 *BioMedical Engineering OnLine*, 22(1):117, 2023.
- 597
598 Yanzhou Su, Tianbin Li, Jiyao Liu, Chenglong Ma, Junzhi Ning, Cheng Tang, Sibo Ju, Jin Ye,
599 Pengcheng Chen, Ming Hu, et al. Gmai-vl-r1: Harnessing reinforcement learning for multimodal
600 medical reasoning. *arXiv preprint arXiv:2504.01886*, 2025.
- 601
602 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqiang Chang, Kai Sheng,
603 Wei Liu, Junyang Wang, et al. Qwen2.5-vl: Enhancing vision-language model's perception of the
604 world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024a.
- 605
606 Qian Wang, Baiqiao Liu, Hongjian Wang, Jiaheng Li, Qingyu Chen, and Zhiyong Lu. Lingshu:
607 A linguistically-enhanced multi-modal chinese medical large language model. *arXiv preprint
arXiv:2406.06489*, 2024b.
- 608
609 Yinhuai Wang, Jiwen Yu, and Jian Zhang. Zero-shot image restoration using denoising diffusion
610 null-space model. *arXiv preprint arXiv:2212.00490*, 2022.
- 611
612 Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment:
613 from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–
612, 2004.
- 614
615 Haoning Wu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Annan Wang, Chunyi Li,
616 Wenxiu Sun, Qiong Yan, Guangtao Zhai, et al. Q-instruct: Improving low-level visual abilities for
617 multi-modality foundation models. In *Proceedings of the IEEE Conference on Computer Vision
and Pattern Recognition*, pp. 5499–5509, 2024a.
- 618
619 Haoning Wu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Annan Wang, Chunyi Li,
620 Wenxiu Sun, Qiong Yan, Guangtao Zhai, et al. Q-bench: A benchmark for general-purpose
621 foundation models on low-level vision. In *ICLR*, 2024b.
- 622
623 xAI Team. Grok-4: Large language model. *Technical Report*, 2024. Available at <https://x.ai>.
- 624
625 Huihui Xu, Yuanpeng Nie, Hualiang Wang, Ying Chen, Wei Li, Junzhi Ning, Lihao Liu, Hongqiu
626 Wang, Lei Zhu, Jiyao Liu, et al. Medground-r1: Advancing medical image grounding via spatial-
627 semantic rewarded group relative policy optimization. *arXiv preprint arXiv:2507.02994*, 2025.
- 628
629 Siyi Xun, Yue Sun, Jingkun Chen, Zitong Yu, Tong Tong, Xiaohong Liu, Mingxiang Wu, and Tao
630 Tan. Mediqa: A scalable foundation model for prompt-driven medical image quality assessment.
631 *arXiv preprint arXiv:2507.19004*, 2025.
- 632
633 Jin Ye, Junlong Cheng, Jianpin Chen, Zhongying Deng, Tianbin Li, Haoyu Wang, Yanzhou Su,
634 Ziyan Huang, Jilong Chen, Lei Jiang, et al. Sa-med2d-20m dataset: Segment anything in 2d
635 medical imaging with 20 million masks. *arXiv preprint arXiv:2311.11969*, 2023.
- 636
637 Jin Ye, Guoan Wang, Yanjun Li, Zhongying Deng, Wei Li, Tianbin Li, Haodong Duan, Ziyan Huang,
638 Yanzhou Su, Benyou Wang, et al. Gmai-mmbench: A comprehensive multimodal evaluation
639 benchmark towards general medical ai. *Advances in Neural Information Processing Systems*, 37:
94327–94427, 2024.
- 640
641 Zhiyuan You, Jinjin Tang, Yao Zhao, Xiaoming Wu, Zhifeng Duanmu, and Zhou Wang. Depict-
642 ing beyond scores: Advancing image quality assessment through multi-modal language models.
643 *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(5):4209–4217, 2024.
- 644
645 Jure Zbontar, Florian Knoll, Anuroop Sriram, Tullie Murrell, Zhengnan Huang, Matthew J Muckley,
646 Aaron Defazio, Ruben Stern, Patricia Johnson, Mary Bruno, et al. fastmri: An open dataset and
647 benchmarks for accelerated mri. *arXiv preprint arXiv:1811.08839*, 2018.
- 648
649 Lin Zhang, Lei Zhang, and Alan C Bovik. A feature-enriched completely blind image quality
650 evaluator. *IEEE Transactions on Image Processing*, 24(8):2579–2591, 2015.

648 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable
649 effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on*
650 *computer vision and pattern recognition*, pp. 586–595, 2018.

651
652 Yizhou Zhang, Tianhe Chen, Shuai Liu, Rui Zhao, and Xiaodong Wang. Bias in artificial intelli-
653 gence for medical imaging: fundamentals, detection, avoidance, mitigation, challenges, ethics,
654 and prospects. *Diagnostic and Interventional Radiology*, 30(3):242854, 2024.

655 Ádám Nárai, Petra Hermann, Tibor Auer, Péter Kemenczky, János Szalma, István Homolya, Eszter
656 Somogyi, Pál Vakli, Béla Weiss, and Zoltán Vidnyánszky. "movement-related artefacts (mr-art)
657 dataset", 2022.

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702	A APPENDIX	
703		
704	APPENDIX TABLE OF CONTENTS	
705		
706	A.1 Data Construction Pipeline and Quality Control	15
707	A.2 Detailed Benchmark Statistics	17
708	A.2.1 Dataset Composition by Source and Modality	17
709	A.2.2 Distributions of Task Types and Degradation Levels in MedQ-Perception	17
710	A.2.3 Distribution of low-level attributions	18
711	A.3 Evaluation Prompt	18
712	A.4 Complete Experimental Results	21
713	A.4.1 Experimental Setup	21
714	A.4.2 Detailed Model Performance	21
715	A.5 Qualitative Analysis and Case Studies	23
716	A.5.1 Why Do Medical-Specialized Models Underperform General-Purpose Models?	23
717	A.5.2 Example of Reasoning Tasks	23
718	A.5.3 Human Expert Evaluation Protocol	26
719		
720		
721		
722		
723		
724		
725		
726		
727		
728		
729		
730		
731		
732		
733		
734		
735		
736		
737		
738		
739		
740		
741		
742		
743		
744		
745		
746		
747		
748		
749		
750		
751		
752		
753		
754		
755		

756 A.1 DATA CONSTRUCTION PIPELINE AND QUALITY CONTROL
757

758 The construction of MedQ-Bench involved a systematic multi-stage pipeline for collecting, curating,
759 and annotating medical images across five modalities. This section provides detailed information
760 about our comprehensive data sourcing strategies, quality control measures, and annotation proto-
761 cols, with particular emphasis on the diverse sources and label types that enable robust evaluation of
762 low-level visual perception capabilities.

763
764 **Comprehensive Data Sources and Acquisition Strategy.** We employed a three-channel data
765 collection strategy: "public datasets + imaging department collaboration + synthetic generation".
766 On one hand, we conducted comprehensive internet searches for 2D/3D medical quality-related
767 datasets; on the other hand, we collaborated with hospitals to obtain ethically approved clinical
768 data. From this massive data pool, we ultimately selected the datasets shown in Table 3, covering
769 5 medical imaging modalities to ensure universality and clinical relevance of data sources. For im-
770 ages, we adhere to the SA-Med2D-20M (Ye et al., 2023) protocol, transforming all 2D/3D medical
771 images into 2D RGB images for further evaluation. Table 3 provides a complete overview of all
772 datasets integrated into MedQ-Bench, including specific modalities, sample quantities, label types,
773 and acquisition status. The table demonstrates the comprehensive scope of our data collection ef-
774 fort, spanning established clinical research datasets and custom synthetic degradation collections,
775 and AI-generated images. All collected images were anonymized, with all patient-identifying infor-
776 mation systematically removed using automated de-identification pipelines validated against clinical
777 privacy requirements.

778
779 **Expert-designed Seed Perception Questions.** The construction process began with a panel of
780 medical imaging specialists who designed seed questions covering diverse modalities, degradation
781 types, and task formats. These domain experts systematically identified key visual quality attributes
782 specific to each modality, ensuring that the seed questions span the three question types: Yes-No,
783 What, and How. Each seed question was carefully paired with selected images to ensure strong
784 alignment between the textual prompt and visual evidence, establishing a foundation of clinically
785 grounded quality assessment scenarios.

786
787 **Controlled Question Expansion.** To scale beyond the initial seed set while maintaining quality
788 and clinical relevance, we employed GPT-4o as a controlled question generator. For each image,
789 using seed questions as templates, we randomly selected one question from each question type
790 and performed controlled generation. This systematic generation process varied degradation types,
791 severity levels, and phrasing styles while preserving clinical realism and explicitly avoiding high-
792 level diagnostic reasoning. The expansion process was constrained by predefined templates and
793 modality-specific quality attributes to ensure consistency and prevent drift away from the intended
794 low-level visual assessment focus.

795
796 **Multi-round Expert Validation.** We manually annotated the answers for the generated questions
797 to ensure correctness, consistency, and alignment with the intended low-level quality assessment
798 labels. All question-answer pairs underwent rigorous multi-stage human annotation and verification
799 using a structured annotation interface, as shown in Figure 7 and 8. The multi-round validation
800 process involved multiple phases of annotation and proofreading:(1) Initial independent review by
801 at least three medical imaging experts for question formulation, answer correctness, and image-
802 question alignment; (2) Cross-validation and proofreading sessions to identify and resolve inconsis-
803 tencies; (3) Final consensus rounds where disagreements were resolved through discussion until
804 unanimous agreement was reached. Finally, the dataset was randomly partitioned by image into
805 development and test sets of equal size.

806
807 **Reasoning Annotation Standards and Workflow.** For the MedQ-Reasoning tasks, we estab-
808 lished specific annotation standards to ensure consistent and clinically relevant quality assessment
809 descriptions. Expert annotators followed a structured reasoning workflow that emphasized system-
810 atic analysis and transparent decision-making processes. The reasoning annotation protocol involved
811 a sequential four-step process: (1) Visual Analysis Phase: Systematic examination of perceptual at-
812 tributes such as noise, blur, artifacts, contrast, and resolution, avoiding any high-level diagnostic in-
813 terpretation; (2) Modality-Specific Assessment: Targeted evaluation of quality dimensions specific

810 to each imaging modality (e.g., streak artifacts in CT, motion artifacts in MRI, staining uniformity
 811 in histopathology), following standardized checklists for each modality type; (3) Quality Classifi-
 812 cation: Application of a three-tier system based on accumulated evidence from steps 1-2: "good"
 813 (no significant quality issues affecting clinical utility), "usable" (minor quality issues that do not
 814 compromise diagnostic accuracy), and "reject" (severe quality degradation requiring repeat imag-
 815 ing); (4) Structured Description Generation: Creation of comprehensive yet concise descriptions
 816 (3-5 sentences) that logically connect the observed visual attributes to the final quality judgment,
 817 ensuring clear reasoning traceability from observation to conclusion. This step-by-step reasoning
 818 flow ensures that all quality assessments follow a consistent analytical framework, with each con-
 819 clusion being explicitly grounded in observable visual evidence rather than subjective impressions.
 820 All reasoning annotations underwent the same multi-round validation process as the perception tasks
 821 to ensure consistency and clinical accuracy across all expert annotators.
 822

823 **Dataset Composition and Balance.** Each modality contributes proportionally to maintain repre-
 824 sentational balance, and degradation types are systematically distributed to avoid bias toward any
 825 particular quality issue.
 826

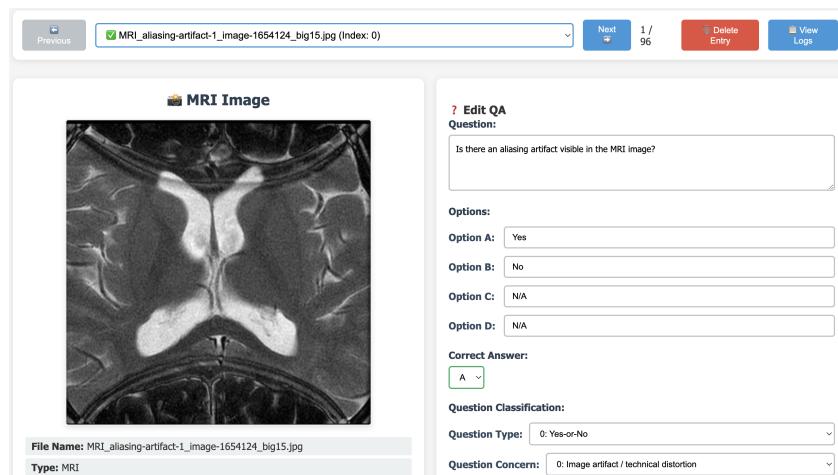


Figure 7: Interface for the MedQ-MCQA dataset.

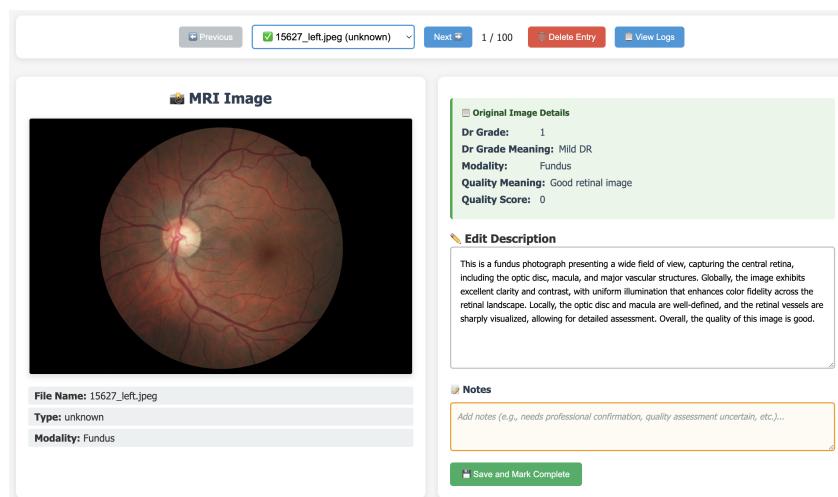


Figure 8: Interface for the MedQ-Reasoning dataset.

864 A.2 DETAILED BENCHMARK STATISTICS
 865

866 A.2.1 DATASET COMPOSITION BY SOURCE AND MODALITY
 867

868 **Dataset Composition.** The MedQ-Bench dataset consists of 3,308 samples distributed across
 869 three primary source types (Table 2). The dataset covers five major medical imaging modalities
 870 with detailed breakdown by source type and specific datasets shown in Table 3.
 871

Source Type	Authentic	Simulate	AI-Generated
Percentage	41.3%	33.9%	24.8%

872 873 874 875 876 Table 2: Distribution of MedQ-Bench dataset by source type.
 877

Modality	Source Type	Dataset	Samples	% of Modality	Total Samples
CT	Authentic	Radioopaedia	239	27.2%	878
	Simulate	AAPM CT-MAR (AAPM CT-MAR Challenge Organizers, 2023)	613	69.8%	
	AI-generated	Alsynthesis (subset)	26	3.0%	
MRI	Authentic	Radioopaedia	130	15.0%	848
	Authentic	MR-ART (Ádám Nárai et al., 2022)	68	7.9%	
	Authentic	FSL Example MRI Artifacts	43	5.0%	
	Simulate	FastMRI Zbontar et al. (2018)	454	54.4%	
	Simulate	5T MRI Data	55	6.4%	
Histopathology	AI-generated	Alsynthesis	98	11.3%	758
	Authentic	HistoArtifacts (Kanwal, 2024)	220	29.0%	
	AI-generated	HARP (Fuchs et al., 2024)	470	62.0%	
Endoscopy	AI-generated	Alsynthesis	68	9.0%	555
	Authentic	EndoCV2020 (Polat et al., 2020)	470	84.7%	
Retinal	AI-generated	Alsynthesis	85	15.3%	269
	Authentic	EyeQ (Fu et al., 2019)	197	73.2%	
Overall Total	AI-generated	Alsynthesis	72	26.8%	269
			3,308		

892 893 894 895 896 Table 3: Comprehensive breakdown of dataset composition.
 897

898 **Detailed Simulation Methods for Synthetic Degradations.** The simulated CT degradations in
 899 AAPM CT-MAR were reconstructed using several algorithms: SIRT, FBP (Kak & Slaney, 2001),
 900 and FISTA (Beck & Teboulle, 2009). Specifically, CT artifacts were systematically simulated to
 901 include three primary degradation types: (1) limited-angle artifacts, (2) metal artifact reduction, and
 902 (3) sparse-view artifacts. For MRI degradations, we primarily simulated acceleration artifacts and
 903 motion artifacts using established computational frameworks. Acceleration artifacts were generated
 904 from private clinical collections using the uMR Jupiter 5T system, obtained under institutional ethi-
 905 cal approval with comprehensive patient anonymization protocols.

906 To generate synthetic images across diverse medical imaging modalities, we employed BAGEL fine-
 907 tuned on domain-specific medical datasets (Deng et al., 2025). This approach ensured that synthetic
 908 degradations maintained clinical realism while providing controlled quality variations essential for
 909 comprehensive benchmark evaluation.

910 A.2.2 DISTRIBUTIONS OF TASK TYPES AND DEGRADATION LEVELS IN MEDQ-PERCEPTION
 911

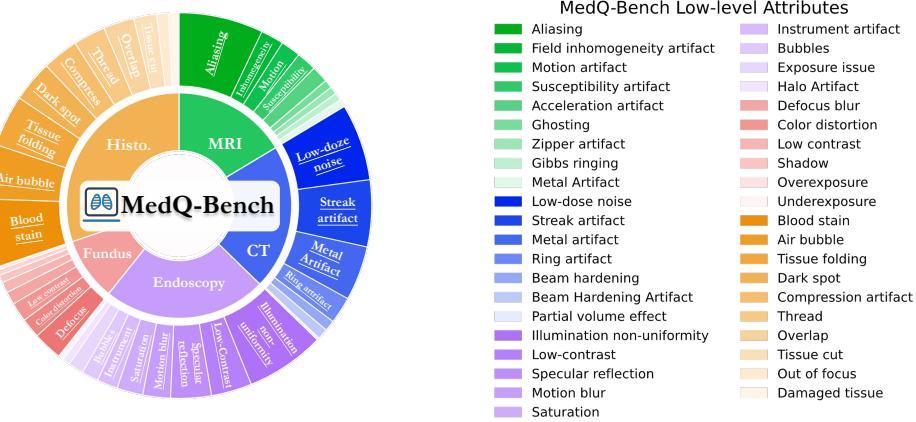
Question Type	Percentage
Modality-specific	57.2%
General	42.8%
Total	100.0%

912 913 914 915 916 917 Table 4: MedQ-Perception: Distribution of
 918 tasks by question type.

Degradation Level	Percentage
No Degradation	23.8%
Mild Degradation	44.6%
Severe Degradation	31.6%
Total	100.0%

919 920 921 922 923 924 Table 5: MedQ-Perception: Distribution of
 925 degradation severity levels.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932



933 Figure 9: Distribution of low-level attributions across imaging modalities and distortion types in
934 MedQ-Bench.

937 A.2.3 DISTRIBUTION OF LOW-LEVEL ATTRIBUTIONS

939 A.3 EVALUATION PROMPT

941 Single-Image Perception Task Prompts

943 Yes-No / What / How Question Template:

945 You are an expert in medical image quality assessment.
946 Please carefully observe this medical image and answer
947 the following question:

950 Reasoning Task Prompts

952 No-reference Reasoning Template:

954 As a medical image quality assessment expert, provide
955 a concise description focusing on low-level appearance
956 of the image in details. Conclude with "Overall, the
957 quality of this image is [good/usable/reject]". Please
958 provide a comprehensive but concise assessment in 3-5
959 sentences.

961 Comprehensive Reasoning Template:

963 As a medical image quality assessment expert, provide
964 a concise description comparing two images focusing
965 on low-level appearance. Conclude with which image
966 has higher quality. Please provide comprehensive but
967 concise assessment in 3-5 sentences.

971 ¹<https://sigpy.readthedocs.io/en/latest/>

²<https://github.com/TorchIO-project/torchio>

972
973

Complete Evaluation Prompt Templates for No-reference Reasoning Tasks

974

Completeness Evaluation Prompt:

975

#System: You are a helpful assistant.
#User: Evaluate whether the description [MLLM DESC] completely includes the low-level visual information in the reference description [GOLDEN DESC]. Please rate score 2 for completely or almost completely including reference information, 0 for not including at all, 1 for including part of the information or similar description.

982

Please only provide the result in the following format: Score:

983

Preciseness Evaluation Prompt:

984

#System: You are a helpful assistant.
#User: The precision metric evaluates whether the low-level description is consistent with the reference and reasonably aligned with the final quality judgment. Minor wording differences or small omissions that do not change the overall meaning should still be considered consistent.
Only penalize clear contradictions with the reference, such as describing blur for clear, noisy for clean, motion-free for motion artifacts, noise-free for low-dose noise, etc.
Evaluate whether output [MLLM DESC] reasonably reflects reference [GOLDEN DESC].
Please rate score 2 for overall consistency and no major contradictions with the quality conclusion, 1 for partial consistency or very few minor contradictions, and 0 for obvious contradictions or misalignment with the quality conclusion.

998

Please only provide the result in the following format: Score:

999

Consistency Evaluation Prompt:

1000

#System: You are a helpful assistant.
#User: Evaluate the internal consistency between the reasoning path (description of image problems) and the final quality judgment in [MLLM DESC]. The reasoning should logically support the final quality conclusion. For example, if many serious problems are described, the final quality should be "reject"; if minor problems are described, it should be "usable"; if no or very few problems are described, it should be "good".
Compare with the reference [GOLDEN DESC] to understand the expected reasoning-conclusion relationship.
Please rate score 2 for highly consistent reasoning and conclusion, 1 for partially consistent with minor logical gaps, and 0 for major inconsistency between described problems and quality judgment.

1014

Please only provide the result in the following format: Score:

1015

Quality Accuracy Evaluation Prompt:

1016

#System: You are a helpful assistant.
#User: Evaluate the accuracy of the final quality judgment in [MLLM DESC] compared to the reference [GOLDEN DESC].

1019

The quality levels have a progressive relationship: reject < usable < good. Consider the distance between predicted and reference quality:

1021

Please rate score 2 for exactly matching the reference quality level, 1 for adjacent level difference (e.g., usable vs good, or reject vs usable), and 0 for distant level difference (reject vs good) or completely incorrect quality assessment.

1024

Please only provide the result in the following format: Score:

1026
1027 Complete Evaluation Prompt Templates for Comparison Reasoning Tasks

1028 **Completeness Evaluation Prompt:**

1029 #System: You are a helpful assistant.

1030 #User: Evaluate whether the description [MLLM DESC] completely

1031 includes the low-level visual information in the reference

1032 description [GOLDEN DESC]. Please rate score 2 for completely

1033 or almost completely including reference information, 0 for not

1034 including at all, 1 for including part of the information or

1035 similar description.

1036 Please only provide the result in the following format: Score:

1037 **Preciseness Evaluation Prompt:**

1038 #System: You are a helpful assistant.

1039 #User: The precision metric evaluates whether the low-level

1040 description is consistent with the reference and reasonably

1041 aligned with the final quality judgment. Minor wording

1042 differences or small omissions that do not change the overall

1043 meaning should still be considered consistent.

1044 Only penalize clear contradictions with the reference, such as

1045 describing blur for clear, noisy for clean, motion-free for

1046 motion artifacts, noise-free for low-dose noise, etc.

1047 Evaluate whether output [MLLM DESC] reasonably reflects

1048 reference [GOLDEN DESC].

1049 Please rate score 2 for overall consistency and no major

1050 contradictions with the quality conclusion, 1 for partial

1051 consistency or very few minor contradictions, and 0 for obvious

1052 contradictions or misalignment with the quality conclusion.

1053 Please only provide the result in the following format: Score:

1054 **Consistency Evaluation Prompt:**

1055 #System: You are a helpful assistant.

1056 #User: Evaluate the internal consistency between the reasoning

1057 path (comparative description of image problems) and the final

1058 quality comparison judgment in [MLLM DESC]. The reasoning

1059 should logically support the final comparison conclusion.

1060 Compare with the reference [GOLDEN DESC] to understand the

1061 expected reasoning-conclusion relationship for image comparison.

1062 Please rate score 2 for highly consistent reasoning and

1063 comparison conclusion, 1 for partially consistent with minor

1064 logical gaps, and 0 for major inconsistency between described

1065 comparative problems and quality comparison judgment.

1066 Please only provide the result in the following format: Score:

1067 **Quality Accuracy Evaluation Prompt:**

1068 #System: You are a helpful assistant.

1069 #User: Evaluate the accuracy of the final quality comparison

1070 judgment in [MLLM DESC] compared to the reference [GOLDEN DESC].

1071 The comparison should correctly identify which image has higher

1072 quality based on the described visual characteristics.

1073 Please rate score 2 for exactly matching the reference quality

1074 comparison, and 0 for completely incorrect quality comparison

1075 (opposite conclusion) or unreasonable assessment.

1076 Please only provide the result in the following format: Score:

1077

1078

1079

1080 **A.4 COMPLETE EXPERIMENTAL RESULTS**

1081 **A.4.1 EXPERIMENTAL SETUP**

1084 In this study, we evaluated various large vision-language models (LVLMs), encompassing medical-
 1085 specialized models, open-source models, and closed-source API general-purpose models. Model
 1086 weights were obtained from their respective official Hugging Face repositories. The evaluation
 1087 work was conducted using the VLMEvalKit framework³.

1088 The evaluation was performed under a "zero-shot" setting. Specifically, our evaluation prompts
 1089 contained no example demonstrations, and models had to complete task reasoning without any related
 1090 training or examples. This approach better tests the models' generalization capabilities and under-
 1091 standing abilities, examining their performance when faced with novel problems. All tests were
 1092 executed on NVIDIA A100 GPUs with 80GB memory.

1093 **A.4.2 DETAILED MODEL PERFORMANCE**

Model (variant)	Perception (Dev)				Reasoning (Dev)				
	Yes-or-No↑	What↑	How↑	Overall↑	Comp.↑	Prec.↑	Cons.↑	Qual.↑	Overall↑
Qwen2.5-VL-Instruct (7B)	62.71%	45.26%	53.93%	56.32%	0.688	0.615	1.869	1.122	4.294
Qwen2.5-VL-Instruct (32B)	64.43%	44.40%	<u>56.20%</u>	57.78%	<u>1.036</u>	0.896	<u>1.959</u>	1.253	<u>5.144</u>
Qwen2.5-VL-Instruct (72B)	74.57%	38.36%	55.37%	60.94%	0.864	0.851	1.860	1.348	4.923
InternVL3 (8B)	<u>75.09%</u>	46.98%	50.62%	60.94%	0.937	0.878	1.864	1.339	5.018
InternVL3 (38B)	70.62%	47.84%	51.24%	59.32%	0.928	0.900	<u>1.900</u>	1.367	5.095
BiMediX2 (8B)	47.77%	28.02%	29.13%	37.29%	0.367	0.376	0.348	0.683	1.774
MedGemma (27B)	62.71%	44.40%	49.59%	54.55%	0.742	0.466	1.652	1.249	4.109
Lingshu (32B)	48.80%	50.86%	53.31%	50.85%	0.629	0.733	1.964	1.059	4.385
Mistral-Medium-3	65.46%	46.12%	52.89%	57.32%	0.937	0.805	1.652	1.389	4.783
Claude-4-Sonnet	67.53%	39.66%	53.93%	57.47%	0.837	0.674	1.810	1.385	4.706
Gemini-2.5-Pro	70.10%	<u>52.16%</u>	46.90%	58.24%	0.810	0.769	1.579	<u>1.548</u>	4.706
GPT-4o	73.54%	48.71%	52.89%	<u>61.40%</u>	0.923	<u>0.936</u>	1.809	1.389	5.057
Grok-4	76.98%	46.55%	63.22%	<u>66.41%</u>	<u>1.036</u>	<u>0.937</u>	1.751	1.484	<u>5.208</u>
GPT-5	78.52%	57.33%	<u>56.61%</u>	66.56%	1.176	1.090	1.756	1.566	5.588

1109 Table 6: Performance of different models on perception and no-reference reasoning tasks (Dev Set).

Model	CT	Histo.	MRI	Endos.	Retinal
GPT-5	71.47%	65.43%	75.90%	60.89%	70.09%
GPT-4o	72.85%	58.33%	64.75%	60.44%	66.67%
Grok-4	70.14%	59.37%	65.93%	64.49%	61.40%
Gemini-2.5-Pro	67.04%	53.40%	60.79%	60.89%	59.83%
Mistral-Medium-3	65.93%	38.58%	61.51%	65.33%	61.54%
Claude-4-Sonnet	64.27%	54.63%	55.04%	65.33%	65.81%
Qwen2.5-VL-72B	65.65%	47.53%	<u>74.82%</u>	66.22%	64.96%
InternVL3-38B	<u>68.14%</u>	48.46%	62.95%	60.44%	70.09%
InternVL3-8B	60.66%	51.54%	61.87%	<u>67.56%</u>	63.25%
Qwen2.5-VL-32B	59.00%	46.30%	66.55%	<u>67.56%</u>	63.25%
Qwen2.5-VL-7B	56.79%	35.49%	65.47%	60.44%	64.96%
MedGemma-27B	66.57%	46.60%	57.55%	56.00%	59.83%
Lingshu-32B	57.89%	35.19%	61.15%	50.67%	48.72%
BiMediX2-8B	41.99%	23.38%	44.36%	38.74%	47.62%

1127 Table 7: Detailed perception accuracy results across five imaging modalities on the test set.

1129 Table 9 provides the complete numerical breakdown of model performance across different com-
 1130 parison difficulty levels and evaluation dimensions corresponding to Figure 6. This comprehensive
 1131 analysis reveals significant performance variations between coarse-grained and fine-grained com-
 1132 parison tasks across all models.

1133 ³<https://github.com/open-compass/VLMEvalKit>

Model	Comp.	Prec.	Cons.	Qual.	Overall
GPT-5	1.376	1.504	1.895	1.609	6.384
GPT-4o	<u>1.113</u>	<u>1.489</u>	1.947	<u>1.669</u>	<u>6.218</u>
Grok-4	<u>1.203</u>	1.203	1.865	1.421	5.692
Gemini-2.5-Pro	1.008	1.180	1.895	1.489	5.572
Mistral-Medium-3	0.932	1.263	1.789	1.414	5.398
Claude-4-Sonnet	0.827	0.992	<u>1.917</u>	1.338	5.074
Qwen2.5-VL-72B-Instruct	0.947	1.158	1.481	1.376	4.962
InternVL3-38B	1.090	1.090	1.684	1.489	5.353
InternVL3-8B	1.023	<u>1.278</u>	<u>1.910</u>	<u>1.549</u>	<u>5.760</u>
Qwen2.5-VL-32B-Instruct	0.865	0.872	1.887	1.083	4.707
Qwen2.5-VL-7B-Instruct	0.684	0.925	1.316	1.150	4.075
MedGemma-27B	0.662	0.571	1.105	0.955	3.293
Lingshu-32B	0.692	0.940	1.519	1.203	4.354
BiMediX2-8B	0.526	0.579	0.594	0.511	2.210

Table 8: Performance comparison on MedQ-Reasoning paired comparison tasks (Dev Set).

Model	Group	Comp.	Prec.	Cons.	Qual.	Acc.	Total
GPT-5	Overall	1.293	1.556	1.925	1.564	6.338	
	Fine-grained	1.301	1.495	1.903	1.476	6.175	
	Coarse-grained	1.267	1.767	2.000	1.867	6.900	
GPT-4o	Overall	1.105	1.414	1.632	1.564	5.714	
	Fine-grained	1.155	1.388	1.583	1.534	5.660	
	Coarse-grained	0.933	1.500	1.800	1.667	5.900	
Grok-4	Overall	1.150	1.233	1.820	1.459	5.662	
	Fine-grained	1.214	1.350	1.825	1.495	5.883	
	Coarse-grained	0.933	0.833	1.800	1.333	4.900	
Gemini-2.5-Pro	Overall	1.053	1.233	1.774	1.534	5.594	
	Fine-grained	1.039	1.262	1.709	1.476	5.485	
	Coarse-grained	1.100	1.133	2.000	1.733	5.967	
InternVL3-8B	Overall	0.985	1.278	1.797	1.474	5.534	
	Fine-grained	0.971	1.155	1.748	1.359	5.233	
	Coarse-grained	1.033	1.700	1.967	1.867	6.567	
Claude-4-Sonnet	Overall	0.857	1.083	1.910	1.481	5.331	
	Fine-grained	0.835	1.107	1.883	1.495	5.320	
	Coarse-grained	0.933	1.000	2.000	1.433	5.367	
Mistral-Medium-3	Overall	0.872	1.203	1.827	1.338	5.241	
	Fine-grained	0.893	1.252	1.786	1.282	5.214	
	Coarse-grained	0.800	1.033	1.967	1.533	5.333	
InternVL3-38B	Overall	1.075	1.083	1.571	1.414	5.143	
	Fine-grained	1.117	1.184	1.466	1.359	5.126	
	Coarse-grained	0.933	0.733	1.933	1.600	5.200	
Lingshu-32B	Overall	0.729	1.015	1.586	1.323	4.654	
	Fine-grained	0.699	0.990	1.505	1.243	4.437	
	Coarse-grained	0.833	1.100	1.867	1.600	5.400	
Qwen2.5-VL-32B	Overall	0.692	0.752	1.895	0.962	4.301	
	Fine-grained	0.786	0.922	1.864	1.068	4.641	
	Coarse-grained	0.367	0.167	2.000	0.600	3.133	
Qwen2.5-VL-7B	Overall	0.714	0.902	1.316	1.143	4.075	
	Fine-grained	0.757	1.000	1.320	1.175	4.252	
	Coarse-grained	0.567	0.567	1.300	1.033	3.467	
Qwen2.5-VL-72B	Overall	0.737	0.977	1.233	1.113	4.060	
	Fine-grained	0.699	0.903	1.029	0.971	3.602	
	Coarse-grained	0.867	1.233	1.933	1.600	5.633	
MedGemma-27B	Overall	0.684	0.692	1.128	1.000	3.504	
	Fine-grained	0.650	0.641	0.942	0.854	3.087	
	Coarse-grained	0.800	0.867	1.767	1.500	4.933	
BiMediX2-8B	Overall	0.474	0.549	0.639	0.511	2.173	
	Fine-grained	0.359	0.379	0.641	0.311	1.689	
	Coarse-grained	0.867	1.133	0.633	1.200	3.833	

Table 9: Detailed numerical results for paired comparison reasoning tasks across models, corresponding to Figure 6.

1188 A.5 QUALITATIVE ANALYSIS AND CASE STUDIES
1189

1190 A.5.1 WHY DO MEDICAL-SPECIALIZED MODELS UNDERPERFORM GENERAL-PURPOSE
1191 MODELS?

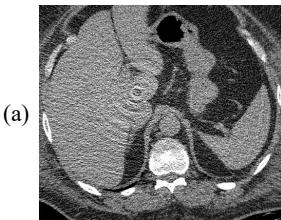
1192 The counterintuitive finding that medical-specialized models consistently underperform general-
1193 purpose models across all evaluation dimensions warrants comprehensive analysis. Figure 11 pro-
1194 vides illustrative examples demonstrating fundamental limitations in medical-specialized models’
1195 low-level visual perception capabilities.
1196

1197 **Insufficient Low-Level Visual Attribute Training.** Medical-specialized models appear to pri-
1198 oritize high-level diagnostic reasoning over fundamental visual perception skills. In the CT scan
1199 example (Figure 11), MedGemma-27B correctly identifies anatomical structures and acknowledges
1200 the presence of streak artifacts, but fails to appropriately assess their clinical significance. The
1201 model describes the image as ”usable but not optimal” despite prominent metal artifacts that would
1202 necessitate repeat scanning in clinical practice. This suggests that medical fine-tuning datasets may
1203 inadequately represent the full spectrum of image quality degradations encountered in clinical work-
1204 flows.

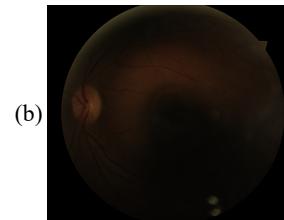
1205 **Diagnostic Bias Over Quality Assessment.** BiMediX2-8B demonstrates a critical failure mode
1206 by describing the same severely degraded CT scan as having ”good quality and suitable for diag-
1207 nosis.” This systematic misalignment indicates that medical-specialized training may inadvertently
1208 optimize models for diagnostic confidence rather than quality assessment accuracy. The model’s
1209 focus on anatomical identification overshadows its ability to detect quality-compromising artifacts,
1210 suggesting that current medical training paradigms may not adequately distinguish between diag-
1211 nostic content recognition and image quality evaluation.

1213 **Question:** Is there a ring artifact visible in the CT
1214 image of the abdomen?
1215 A. Yes, B. No.

1216 **Correct Answer:** A



1213 **Question:** What is the primary quality issue of this image?
1214 A. Blurry and out of focus, B. Poor contrast and tonal separation,
1215 C. Excessive noise and artifacts, D. Insufficient illumination
1216 **Correct Answer:** D



- 1224 ✓ InternVL3-8B: A
1225 ✓ Qwen2.5-VL-32B-Instruct: A.
1226 ✗ Lingshu-32B: B.
1227 ✗ mistral-medium-3_trans: B
1228 ✓ BiMediX2-8B: A
1229 ✓ GPT4o_trans: A
1230 ✓ Claude4_Sonnet_trans: A
1231 ✓ InternVL3-38B: A
1232 ✓ gpt-5: A
1233 ✗ Qwen2.5-VL-72B-Instruct: B
1234 ✗ Qwen2.5-VL-7B-Instruct: B.
1235 ✓ MedGemma-27B: A
1236 ✓ grok-4: A
1237 ✓ GeminiPro2-5_trans: A

- 1224 ✗ InternVL3-8B: B
1225 ✓ Qwen2.5-VL-32B-Instruct: D
1226 ✗ Lingshu-32B: B.
1227 ✗ mistral-medium-3_trans: B. Poor contrast and tonal separation
1228 ✗ BiMediX2-8B: B
1229 ✓ GPT4o_trans: D
1230 ✓ Claude4_Sonnet_trans: D
1231 ✗ InternVL3-38B: A
1232 ✓ gpt-5: D
1233 ✓ Qwen2.5-VL-72B-Instruct: D
1234 ✓ Qwen2.5-VL-7B-Instruct: D
1235 ✓ MedVLM-R1: D
1236 ✓ grok-4: D
1237 ✓ GeminiPro2-5_trans: D

1238 Figure 10: Representative QA examples demonstrating typical question-answer patterns in MedQ-
1239 Bench across different medical imaging modalities and quality assessment scenarios.
1240

1241 A.5.2 EXAMPLE OF REASONING TASKS

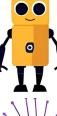
- 1242
 1243
 1244
 1245  **User:** As a medical image quality assessment expert, provide a concise description focusing on low-level appearance of the image in details. Conclude with "Overall, the quality of this image is [good/usable/reject]".
- 1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
- 1254  **Golden Reasoning:** This is a CT scan of the abdominal region. The image exhibits moderate clarity and contrast, but the presence of significant noise reduces overall resolution. Metal artifacts are prominent, creating streaks that obscure anatomical details and hinder accurate interpretation of the abdominal structures. Overall, the quality of this image is reject.
- 1255
 1256
 1257
- 1258  **InternVL3-8B:** This abdominal CT scan shows severe streak artifacts radiating from the center, likely due to beam hardening or metal. While overall contrast is adequate and major structures like vertebrae and soft tissues are visible, the artifacts obscure finer details and reduce clarity. Some noise is present but less problematic than the artifacts. Overall, the image quality is poor and judged as reject.
- 1259
 1260
 1261
- 1262  **MedGemma-27B:** This abdominal CT scan shows bones, soft tissues, and air with adequate brightness and contrast. However, streak artifacts from dense structures, such as bone or metal, obscure fine details. Moderate noise adds a grainy texture, further reducing clarity. While the image remains interpretable for basic anatomical identification, the artifacts and noise limit diagnostic precision. Overall, the quality is usable but not optimal.
- 1263
 1264
 1265
- 1266  **Lingshu-32B:** This CT scan of the abdomen or pelvis shows the vertebral column, soft tissues, and possible lesions. However, it is heavily degraded by bright streak artifacts that obscure much of the anatomy. These distortions suggest poor acquisition parameters or technical issues. As a result, the image is of low quality and not reliable for diagnostic purposes.
- 1267
 1268
 1269
- 1270  **BiMedIX2-8B:** This CT scan of the liver shows clear structures with good contrast and minimal noise. The organ and surrounding tissues are well-defined, supporting accurate anatomical assessment. No significant artifacts are visible, and the overall resolution is sufficient for clinical use. Overall, the image quality is good and suitable for diagnosis.
- 1271
 1272
 1273
- 1274  **GPT-40:** The abdominal CT image contains multiple bright streak artifacts, likely from metal or dense materials. These streaks obscure anatomical details and reduce clarity. Uneven brightness and contrast further complicate interpretation. Due to the severity of these distortions, the overall image quality is reject.
- 1275
 1276
 1277
- 1278  **Claude4 Sonnet:** This CT scan is heavily degraded by extensive streak artifacts radiating across the field of view, creating a star-like pattern. The artifacts, likely from metallic implants, obscure soft tissue structures and organ boundaries. Contrast is poor, noise is elevated, and diagnostic interpretation is nearly impossible in affected regions. While a few landmarks remain visible, the image quality overall is reject.
- 1279
 1280
 1281
- 1282  **InternVL3-38B:** This abdominal CT scan shows vertebrae, soft tissues, and organs with adequate contrast and sharpness. However, background noise is high, which reduces the clarity of smaller structures and subtle abnormalities. No major artifacts obstruct critical regions, and the field of view is appropriate. Overall, the image is usable but limited by noise.
- 1283
 1284
 1285
- 1286  **GPT-5:** This axial CT image shows moderate noise and acceptable sharpness, with organ edges reasonably defined. However, severe metal-induced streak artifacts radiate across the field, producing alternating bright and dark bands that obscure anatomy. Global shading and banding further reduce uniformity. While centered and properly cropped, the distortions make diagnostic interpretation unreliable. Overall, the quality is reject.
- 1287
 1288
 1289
- 1290
 1291
- 1292  **Qwen2.5-VL-72B-Instruct:** The CT scan is degraded by strong streak artifacts, likely from metallic implants or dense objects. These streaks obscure anatomical details and compromise evaluation of affected regions. Contrast and brightness remain acceptable in unaffected areas, preserving partial visibility of tissues. Despite this, the diagnostic utility is poor. Overall, the quality is reject.
- 1293
 1294
 1295

Figure 11: Representative no-reference reasoning image examples demonstrating typical question-answer patterns in MedQ-Bench across different models.

1296

1297

1298



1299

User: As a medical image quality assessment expert, provide a concise description comparing two images focusing on low-level appearance. Conclude with which image has higher quality.

1300

1301

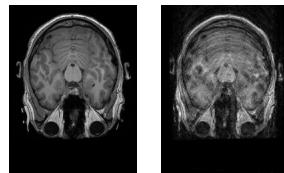
1302

1303

1304

1305

1306



1307

1308

1309

1310

1311



Golden Reasoning: The first image demonstrates excellent clarity and sharpness with minimal motion artifacts, allowing for clear visualization of anatomical structures. In contrast, the second image exhibits severe motion artifacts, including blurring and ghosting, which significantly obscure anatomical details and reduce diagnostic quality. The brain tissue boundaries and ventricular structures are clearly visible in the first image, while they are severely degraded in the second. Overall, the first image has higher quality.

1312

1313

1314



InternVL3-8B: The comparison of the two images shows clear quality differences. Image 1 has better contrast and clarity, with gray and white matter well-differentiated and minimal noise, making structures like gyri and sulci more distinct. Its resolution is higher, and artifacts are minimal, ensuring reliable diagnostic value. Image 2, by contrast, suffers from lower contrast and clarity, higher noise, and poorer resolution, making brain structures harder to interpret. It also contains more artifacts, which further obscure details. Conclusion: Image 1 offers superior diagnostic quality, providing clearer visualization of anatomical structures.

1315

1316

1317

1318

1319



MedGemma-27B: Image 1 shows clearer brain structures with better gray-white matter contrast and less noise. Image 2 appears blurrier, with weaker contrast and more graininess, making details harder to identify. Overall, Image 1 demonstrates superior clarity and quality.

1320

1321

1322

1323



Lingshu-32B: The first image is a coronal T1-weighted MRI with clear anatomical structures, good signal-to-noise ratio, and strong contrast between gray and white matter. These qualities make it suitable for accurate interpretation. The second image, however, is heavily corrupted by noise, which obscures anatomical details and reduces clarity. The poor visibility significantly limits its diagnostic value. Overall, the first image is of much higher quality.

1324

1325

1326



BiMediX2-8B: The left image demonstrates a higher signal-to-noise ratio, producing a clearer and more detailed view of anatomical structures. In contrast, the right image shows lower SNR, resulting in reduced clarity and less visible details. This difference strongly impacts the ability to detect and diagnose subtle abnormalities. Overall, the left image provides superior quality and greater diagnostic reliability.

1327

1328

1329

1330

1331

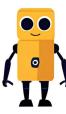


GPT-40: The first image shows higher contrast, clearer anatomical structures, and less noise, with well-defined edges. The second image has lower contrast and more noise, causing structures to appear blurred and less distinct. Overall, the first image demonstrates superior quality due to its clarity and better contrast.

1332

1333

1334



Claude4 Sonnet: The left MRI image has higher contrast, sharper anatomical boundaries, and a better signal-to-noise ratio, with clear visualization of brain parenchyma and ventricles. Its signal is more uniform and details are well-preserved. The right image, however, appears grainier with lower contrast, more noise, and blurred structures, showing reduced diagnostic value. Overall, the left image is of higher quality and more suitable for clinical interpretation.

1335

1336



InternVL3-38B: The first MRI image shows clear brain structures with distinct tissue contrasts and high resolution, making fine details easy to visualize. In comparison, the second image appears lower in quality, with weaker tissue contrast and less distinct details. Its resolution is reduced, making fine anatomical features harder to discern. Overall, the first image demonstrates higher quality due to sharper details and better contrast.

1337

1338

1339

1340



GPT-5: The left image has sharp cortical and orbital boundaries, low noise, and strong gray-white matter contrast. Its background is smooth and uniform, with minimal artifacts. By contrast, the right image is blurred, noisy, and shows reduced contrast with visible mottling and ghosting. These issues obscure fine details and reduce diagnostic value. Overall, the left image is of clearly higher quality.

1341

1342

1343

1344



Qwen2.5-VL-72B-Instruct: Image 1 has lower contrast, less differentiation between tissues, and smoother appearance with reduced sharpness. Image 2 shows higher contrast and clearer anatomical details, though it also contains more visible noise. Despite the noise, the sharper features and stronger tissue contrast make Image 2 more suitable for diagnostic use. Overall, Image 2 is of higher quality.

1345

1346

1347

1348

1349

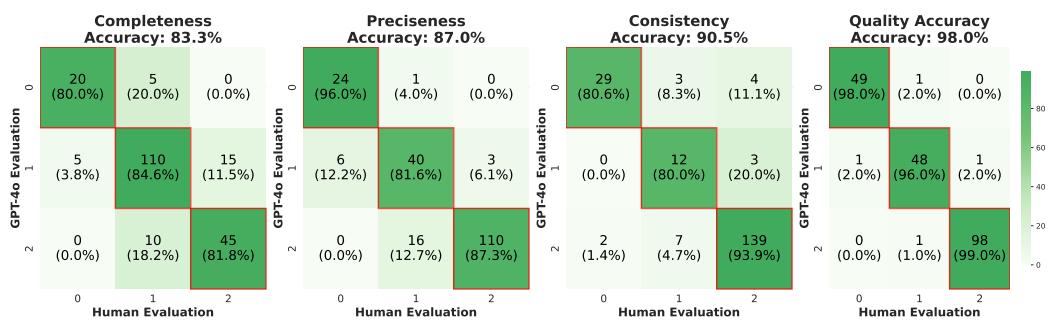
Figure 12: Representative paired image examples demonstrating typical question-answer patterns in MedQ-Bench across different models.

1350 A.5.3 HUMAN EXPERT EVALUATION PROTOCOL

1351
1352 Expert Recruitment and Qualification Criteria. Human experts in our evaluation consisted of
 1353 medical imaging technicians with a minimum of 3 years of clinical experience in medical imaging
 1354 quality assessment and medical imaging PhDs with specialized training in image quality evaluation.
 1355 Medical imaging technicians were recruited from certified clinical facilities and possessed active
 1356 professional certifications in their respective imaging modalities. PhDs were selected from accred-
 1357 ited medical imaging research programs and had completed at least 2 years of coursework, including
 1358 medical image processing and quality assessment methodologies. All experts demonstrated profi-
 1359 ciency in identifying common imaging artifacts and quality issues across multiple medical imaging
 1360 modalities through a standardized pre-evaluation assessment.

1361
1362 Human-AI Alignment Analysis. The confusion matrices shown in Figure 13 demonstrate strong
 1363 alignment between human expert scores and GPT-4o automated evaluation across all three evalua-
 1364 tion dimensions, with over 80% accuracy in each dimension. Quadratic weighted κ_w accounts for
 1365 the ordinal nature of the evaluation labels, penalizing larger discrepancies more heavily than adjacent
 1366 category differences. The consistently high κ_w values (0.774–0.985) detailed in Table 10 indicate
 1367 substantial agreement beyond chance between human expert scores and GPT-4o automated evalua-
 1368 tion, reflecting that the automated system is not only accurate but also aligned with the fine-grained
 1369 ordinal structure of human expert judgments.

1370 The consistently high κ_w values (0.774–0.985) detailed in Table 10 indicate substantial agreement
 1371 beyond chance between human expert scores and GPT-4o automated evaluation, reflecting that the
 1372 automated system is not only accurate but also aligned with the fine-grained ordinal structure of
 1373 human expert judgments. This confirms that our automated evaluation framework maintains robust
 1374 alignment with human expert annotations, strengthening confidence in its use as a reliable surrogate
 1375 for large-scale human evaluation.



1387
 1388 Figure 13: Confusion matrices showing alignment between human expert scores and GPT-4o auto-
 1389 mated evaluation across four evaluation dimensions.

Metric	Completeness	Preciseness	Consistency	Quality Accuracy
κ_w	0.774	0.876	0.840	0.985

1394
 1395 Table 10: Quadratic weighted Cohen’s κ_w values for human–AI alignment across evaluation dimen-
 1396 sions.