

U2-BENCH: Benchmarking Large Vision-Language Models on Ultrasound Understanding

Anjie Le^{*1,2}, Henan Liu^{*1,3}, Yue Wang⁹, Zhenyu Liu¹, Rongkun Zhu⁴,
Taohan Weng^{1,3}, Jinze Yu^{1,3}, Boyang Wang^{1,3}, Yalun Wu³, Kaiwen Yan⁹,
Quanlin Sun⁵, Meirui Jiang^{1,7}, Jialun Pei⁷, Siya Liu¹, Haoyun Zheng¹, Zhoujun Li³,
J. Alison Noble², Jacques Souquet^{1,6}, Xiaoqing Guo^{†4,2}, Manxi Lin^{†8}, Hongcheng Guo^{†1,3}

¹Dolphin AI ²University of Oxford ³Beihang University
⁴Hong Kong Baptist University ⁵University of Cambridge ⁶Chinese Academy of Sciences
⁷The Chinese University of Hong Kong ⁸Technical University of Denmark
⁹Independent

Abstract

Ultrasound is a widely-used imaging modality critical to global healthcare, yet its interpretation remains challenging due to variability in image quality caused by operator dependency, noise, and anatomical complexity. Although large vision-language models (LVLMs) have demonstrated impressive multimodal capabilities across natural and medical domains, their performance on ultrasound remains largely unexplored. We introduce U2-BENCH, the first comprehensive benchmark to evaluate LVLMs on ultrasound understanding across classification, detection, regression, and text generation tasks. U2-BENCH aggregates 7,241 cases spanning 15 anatomical regions and defines 8 clinically inspired tasks, such as *diagnosis*, *view recognition*, *lesion localization*, *clinical value estimation*, and *report generation*, across 50 ultrasound application scenarios. We evaluate 20 state-of-the-art LVLMs, both open- and closed-source, general-purpose and medical-specific. Our results reveal strong performance on image-level classification, but persistent challenges in spatial reasoning and clinical language generation. U2-BENCH establishes a rigorous and unified testbed to assess and accelerate LVLM research in the uniquely multimodal domain of medical ultrasound imaging.¹

1 Introduction

“In diagnostics, the eye sees what the mind knows – true understanding requires merging image patterns with clinical wisdom.” – William Osler

Ultrasound (US) is one of the most widely used imaging modalities in global healthcare — essential in obstetrics, emergency medicine, cardiology, and low-resource settings — while its interpretation remains notoriously difficult [29]. Compared to modalities such as computed tomography (CT), magnetic resonance imaging (MRI), positron emission tomography (PET), and whole-slide imaging (WSI), which offer higher spatial resolution, consistent image quality, and standardized anatomical views, ultrasound is real-time and low-cost but highly operator-dependent and frequently affected by imaging artifacts [65]. In addition, in contrast to these modalities, US is dynamically presenting three-dimensional (3D) anatomies in image sequences. Therefore, accurate interpretation of US demands not only visual pattern recognition in the images, but also an understanding of anatomy and

¹Codes are available at: <https://anonymous.4open.science/r/U2-Bench-F781/VLMEVALKIT/>
Data is available at: <https://huggingface.co/datasets/DolphinAI/u2-bench/tree/main>

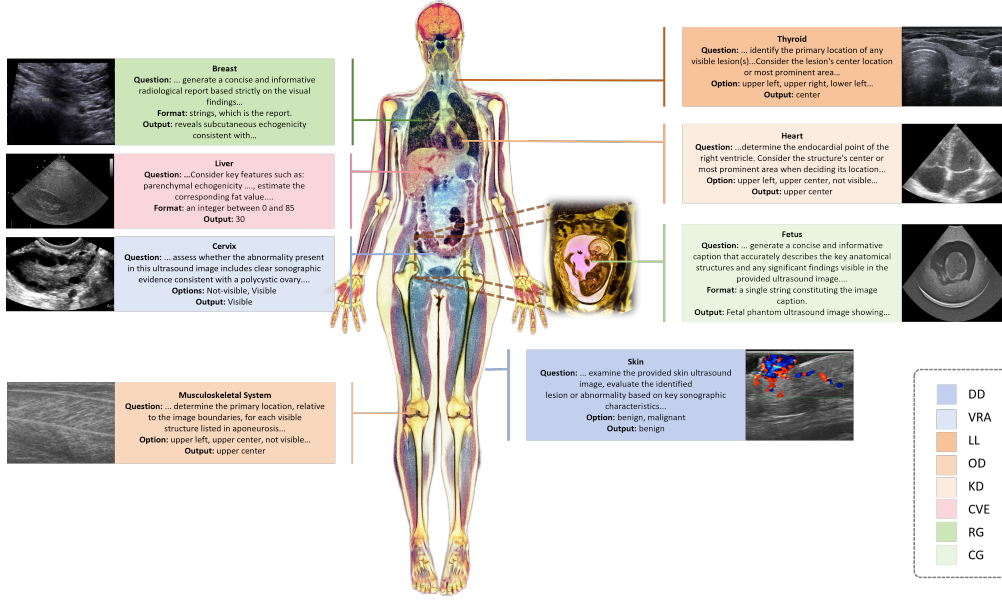


Figure 1: Examples of the 8 benchmark tasks in U2-BENCH across diverse anatomical regions. Each callout, consisting of the question prompt, expected output format, and sample output, highlights a representative ultrasound application scenario of the corresponding task. Tasks involve disease diagnosis (DD), view recognition and assessment (VRA), lesion localization (LL), organ detection (OD), keypoint detection (KD), clinical value estimation (CVE), report generation (RG) and caption generation (CG).

capturing of dynamic spatial-context reasoning, typically requiring extensive prior domain expertise [74]. These challenges have limited the applicability of earlier artificial intelligence (AI) models. However, recent advances in medical large vision-language models (LVLMs) have shown promise in overcoming these barriers [17, 80, 31], potentially offering a robust multimodal understanding of complex, noisy, and context-rich ultrasound data.

While progress in medical LVLM has been rapid, most previous models and benchmarks focus on those less noisy and static imaging modalities [34, 32, 69], leaving the complexities of ultrasound largely unaddressed. Prior efforts in ultrasound AI are typically based on small, task-specific datasets [81], such as fetal plane identification [25] or pathology segmentation [33, 59]. As model capabilities grow, a public, balanced benchmark for ultrasound understanding is needed to evaluate whether emerging LVLMs can generalize beyond static medical vision tasks, to those requiring spatial reasoning and contextual understanding of anatomical structures.

To address these challenges, we introduce **U2-BENCH**, the first benchmark holistically evaluating current LVLMs for ultrasound understanding across diverse tasks and anatomies. The dataset we use comprises 7,241 cases across 15 anatomical regions, involving breast, heart, lung, etc, covering 8 diverse clinical tasks and 50 application scenarios. Each task belongs to one of the four categories: (1) classification (i.e., disease diagnosis, view recognition and assessment), (2) detection (i.e., lesion localization, organ detection, keypoint detection), (3) regression (i.e., clinical value estimation), (4) text generation (i.e., report generation, caption generation). Samples are selected to ensure balance across data sources, anatomies, and task types, to enable robust evaluation and alleviate dataset-specific bias. Several examples in our **U2-BENCH** are shown in Fig. 1.

We benchmark 20 LVLMs, including both open- and closed-source, general-purpose and medical-specialized models, on a diverse set of US tasks. **U2-BENCH** makes the following key contributions:

- **Comprehensive Dataset:** We release the first publicly available benchmark comprising 7,241 ultrasound cases spanning 15 anatomies and 8 clinical tasks, covering 50 application scenarios. Each case is annotated with task-aligned labels in a unified format and paired with carefully designed prompts, enabling standardized and reproducible evaluation.
- **Task Suite and Evaluation:** We define an eight-task taxonomy spanning *disease diagnosis*, *view recognition and assessment*, *lesion localization*, *organ detection*, *keypoint detection*,

clinical value estimation, report generation, and caption generation. Each task reflects real-world clinical workflows and is paired with standard evaluation metrics. We also introduce an aggregate metric to provide a unified assessment of each model’s overall capability in ultrasound understanding.

- **Empirical Insights:** We conduct the first large-scale evaluation of LVLMs on ultrasound, uncovering consistent trends across model families: models achieve strong performance on image-level disease diagnosis and clinical value estimation tasks, but degrade on spatial reasoning tasks such as view recognition and organ detection. Clinical report generation tasks remain particularly challenging. Performance gains from model scaling can be limited, and compact models occasionally outperform larger ones on certain tasks, suggesting that targeted training may be more impactful than scale alone in ultrasound understanding.

2 Related Work

Recent advances in LVLMs have catalyzed new opportunities in medical image understanding. Our work intersects two key directions: (1) the construction of benchmarks for systematically evaluating LVLM capabilities across multimodal medical tasks, and (2) unlocking the potential of LVLMs for ultrasound imaging.

Large Vision-Language Models. LVLMs such as GPT-4V [54], Claude [5], Gemini [4], DeepSeek-VL [21], LLaVA [45], Qwen-VL [9], and MiniGPT4 [89] have emerged as general-purpose multimodal systems capable of handling tasks like image captioning, visual question answering, and multimodal reasoning. These models are trained on large-scale image-text pairs [66, 62], and their performance has been extensively evaluated in domains such as question answering, mathematics, and science [15, 70, 75, 30, 46]. However, their clinical reliability remains underexplored.

To address this gap, several medical-specialized LVLMs have been proposed. MiniGPT-Med [78] focuses on X-ray, CT, and MRI for tasks such as medical report generation, VQA, and disease identification. RadFM [79] further supports both 2D and 3D modalities. MedDr [26] extends to radiology, pathology, dermatology, retinography, and endoscopy, introducing a retrieval-augmented diagnosis strategy. Yet, these models exclude ultrasound. Med-Gemini [71] spans numerous modalities including ultrasound, though its capability in this domain is limited to caption generation.

Multimodal Benchmarks for Large Vision-Language Models. Several benchmarks assess general-domain LVLMs. MMBench [47], MMT-Bench [85], and SEED-Bench [41] evaluate general-domain LVLMs through bilingual multiple-choice questions, large-scale visual reasoning tasks, and generative comprehension across image and video VQA, respectively. However, these benchmarks emphasize general-purpose visual understanding and omit clinically grounded evaluation.

Early medical VQA datasets like VQA-RAD [37], VQA-Med [10], and PathVQA [27] offer radiology or pathology image-question pairs but are not designed for evaluating LVLMs. GMAI-MMBench [16] introduces a large-scale VQA-style benchmark for medical LVLMs, yet it contains only about 1.4k ultrasound cases primarily focused on classification and segmentation on 6 anatomies, and does not evaluate broader model capabilities such as clinical value estimation or structured report generation. In contrast, our **U2-BENCH** focuses exclusively on ultrasound and includes a diverse set of clinically meaningful tasks and anatomical regions.

3 U2-BENCH

Overview. **U2-BENCH** is designed to holistically assess the capabilities of LVLM in ultrasound tasks. Section 3.1 introduces the eight clinically inspired tasks involved in our evaluation, which reflect essential diagnostic and reasoning abilities in ultrasound understanding. Section 3.2 details our benchmark construction pipeline, including dataset curation, preprocessing, and task-specific prompting. Section 3.3 summarizes the statistical property of the resulting dataset, which comprises 7,241 cases across 15 anatomies.

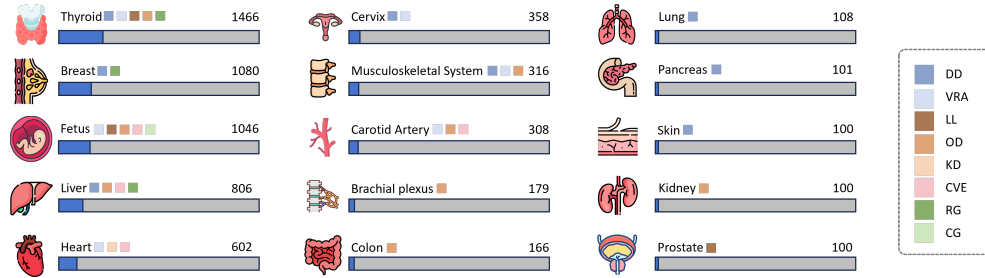


Figure 2: **Distribution of benchmark tasks across 15 anatomical regions in U2-BENCH.** The colored boxes next to each anatomy name indicate the benchmark tasks available for that anatomy, with each color corresponding to one of the eight core tasks (legend shown on the right). The length of the bar reflects the number of samples for each anatomical region. Multiple tasks may share samples from the same anatomical region, depending on annotation availability and clinical relevance.

3.1 Task Definitions

U2-BENCH focuses on four core capabilities: classification, detection, regression, and text generation, to systematically evaluate the performance of LVLMs on ultrasound-related tasks. We define eight tasks based on common ultrasound use cases, designed to probe a range of multimodal abilities, including anatomy recognition and clinical reporting. The task set was informed by typical sonography workflows and refined with input from domain experts to ensure practical relevance. Together, these tasks provide a structured benchmark for assessing LVLM performance across diverse ultrasound application scenarios. The eight tasks are as follows:

Disease Diagnosis (DD). This task requires the model to identify the presence and severity of a disease condition, such as grading in the Breast Imaging Reporting and Data System, based on the appearance of the ultrasound image. The task evaluates the ability of LVLMs to extract high-level semantic features and generate clinically aligned diagnostic predictions.

View Recognition and Assessment (VRA). In clinical practice, accurate diagnosis relies on the clear presentation of anatomical structures from specific angles, referred to as ultrasound standard planes. This task evaluates the ability of a model to assess image quality and classify scans into standard planes corresponding to different anatomical structures, such as the fetal head or abdominal long axis.

Lesion Localization (LL). Given a diagnostic image, the LVLM is asked to identify the location of a lesion, such as a suspicious breast mass, by selecting from nine predefined spatial categories such as upper left, center, or lower right. This task evaluates the spatial reasoning, saliency alignment, and ability to detect subtle structural abnormalities of LVLMs.

Organ Detection (OD). This task involves identifying the presence and boundaries of target organs in the ultrasound field of view, such as liver, kidney, or nerve. It assesses coarse-grained visual recognition under challenges unique to ultrasound, such as acoustic shadowing, inter-patient variability, and orientation ambiguity from manual probe handling.

Keypoint Detection (KD). In measurement tasks such as fetal biometry and adult echocardiography, precise localization of anatomical landmarks is critical for deriving clinically meaningful measurements. This task evaluates the fine-grained spatial understanding and geometric reasoning ability of the model, which are essential for tasks like skeletal length and chamber size estimation.

Clinical Value Estimation (CVE). This task involves predicting continuous clinical parameters derived from ultrasound images, such as lesion size, left ventricular ejection fraction, or liver fat percentage. It covers both anatomical and functional indicators relevant to diagnosis, treatment planning, and longitudinal monitoring, and evaluates whether the model can perform image-to-value regression by mapping visual inputs to clinically meaningful quantitative outputs.

Report Generation (RG). The model is prompted to generate a structured clinical report based on visual input, following the format of example reports provided in the prompt. This task evaluates the ability of LVLM to perform medical language generation and produce outputs that align with standard ultrasound reporting practices.

Caption Generation (CG). The model is asked to generate a concise anatomical description of a diagnostic image, guided by example captions provided in the prompt. This task evaluates basic

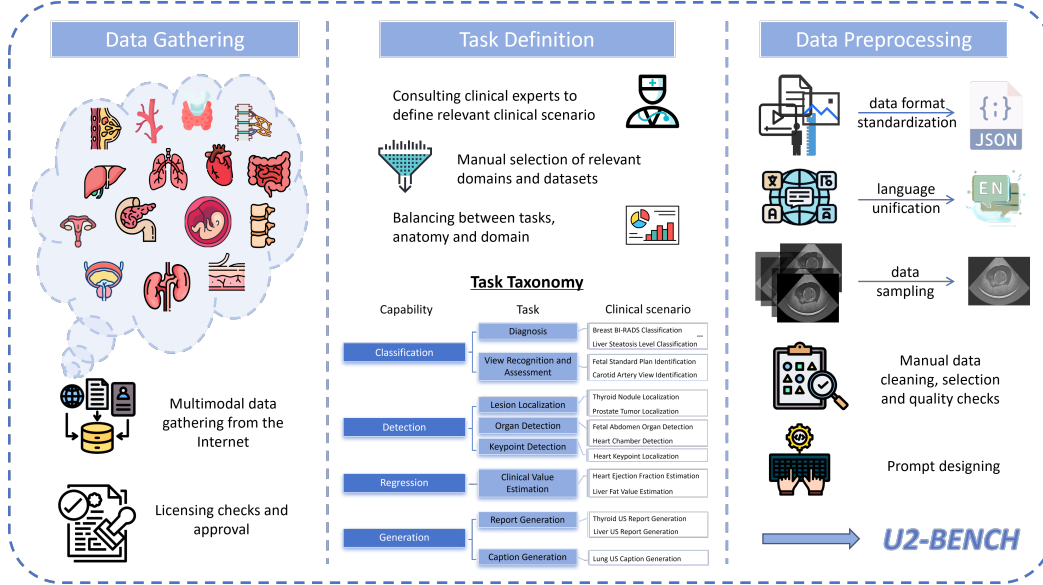


Figure 3: **Overview of the U2-BENCH construction pipeline.** The benchmark is built through three stages: (1) data gathering from 40 licensed ultrasound datasets spanning 15 anatomical regions, (2) task definition across 8 clinically inspired tasks grouped into four core capabilities: classification, detection, regression, and text generation, (3) data preprocessing, including annotation standardization, metadata unification, image/frame selection, and quality verification. This unified pipeline ensures benchmark consistency and clinical relevance across diverse ultrasound scenarios.

visual-language alignment and the ability to verbalize structural features in a clinically appropriate manner of LVLM.

3.2 Data Curation and Processing

In this section, following the approach of previous benchmark constructions [18, 83, 88], we outline the three key steps used to build **U2-BENCH**: (1) data collection and sampling (2) data cleaning, format unification and quality verification, and (3) task-specific prompt design. Figure 3 summarizes the data processing pipeline.

Data Selection and Sampling. We construct **U2-BENCH** by sampling 7,241 ultrasound studies from 40 licensed datasets. These datasets were selected to represent a wide range of diagnostic tasks, anatomical regions, and clinical contexts. While the original datasets were independently curated and clinically annotated, we performed standardization, sampling, and quality checks to ensure consistency across tasks and enable reliable, reproducible benchmarking. Some datasets contribute to multiple benchmark tasks based on their available annotations and clinical relevance.

To reflect real-world clinical data distributions and prevent data leakage, we adopt a task-specific, patient-level sampling strategy. Sampling is performed at the subject level rather than the image level to preserve intra-patient consistency. To ensure anatomical coverage, we include data from 15 anatomical regions: fetus, thyroid, breast, heart, liver, cervix, carotid artery, musculoskeletal system, kidney, prostate, skin, lung, pancreas, brachial plexus, and colon.

Data Cleaning, Format Unification, and Quality Verification. All data in **U2-BENCH** are standardized into a unified format to support consistent parsing and evaluation across the dataset. Ultrasound scans are converted to a uniform image format. For video sequences, a small number of representative frames are sampled per study to control evaluation cost while retaining key diagnostic content. Task-relevant metadata, including anatomy labels, measurements, and reports, is preserved in a structured schema. Segmentation masks are converted to bounding boxes.

To ensure the reliability of **U2-BENCH**, we adopt both automated and manual quality assurance procedures during data preparation.

Table 1: **Summary of annotated datasets used in U2-BENCH, grouped by core capability and task.** The “Case Number” column indicates the number of samples per task, while “Total” reflects the overall count when available. More details about the datasets are included in Appendix F.

| Capability | Task | Case Number | Source Dataset | Total |
|----------------|------|-------------|--|-------|
| Classification | DD | 1,411 | Breast Lesion Detection in Ultrasound Videos [44]; Breast Ultrasound Images Dataset [2]; Dermatologic Ultrasound Images for classification [38]; Knee ultrasound dataset in a population-based cohort [53]; KFGNet [52]; GDPHYSUCC [51]; LEPset [43]; COVID-BLUES [76]; Ultrasound Guided Regional Anesthesia [72]; Ultrasound Breast Images for Breast Cancer [60]; Algerian Ultrasound Images Thyroid Dataset: AUITD [49]; Auto-PCOS classification [22] | 2,999 |
| | VRA | 1,588 | FETAL PLANES DB [12]; FPUS23 [58]; CAMUS [39]; Knee ultrasound dataset in a population-based cohort [53]; Thyroid [36]; ACOUSLIC-AI [61]; JNU-IFM [48]; Carotid Artery Ultrasound and Color Doppler [55]; Auto-PCOS classification [49]; African Fetal Standard Plane [63]; DDTI [57]; CAMUS [39]; CUBS [50]; COVID-BLUES [76]; Dataset of B-mode fatty liver ultrasound images [13]; The Open Kidney Ultrasound Dataset [68]; Micro-Ultrasound Prostate Segmentation Dataset [64]; Breast Ultrasound Images Dataset [2]; Knee ultrasound dataset in a population-based cohort [53]; Polycystic Ovary Ultrasound Images Dataset [77] | |
| Detection | LL | 503 | DDTI [57]; Micro-Ultrasound Prostate Segmentation Dataset [64]; Breast Ultrasound Images Dataset [2]; KFGNet [52]; BrEaST [56] | 2,921 |
| | OD | 1,918 | The Open Kidney Ultrasound Dataset [68]; Echogenic [19]; FALLMUD [24]; CAMUS [39]; HC18 [28]; Thyroid [36]; CCA [11]; Ultrasound Guided Regional Anesthesia [72]; C-TRUS Dataset [40]; ACOUSLIC-AI [61]; PSFHS [7]; JNU-IFM [48]; US simulation & segmentation [73] | |
| | KD | 500 | Unity Imaging Collaborative [67] | |
| Regression | CVE | 521 | CAMUS [39]; CUBS [50]; HC18 [28]; ACOUSLIC-AI [61]; Dataset of B-mode fatty liver ultrasound images [13] | 521 |
| Generation | RG | 600 | Chinese Ultrasound Report Dataset [42] | 800 |
| | CG | 200 | FPUS23 [58] | |
| Overall Total | | | | 7,241 |

(1) Automated Filtering. During data preprocessing, we systematically check for missing labels, inconsistent or invalid annotations, and corrupted or unreadable files. Samples that fail these checks are discarded.

(2) Manual Verification. A team of 10 annotators manually reviewed all cases using a cross-validation protocol, where each data point was independently assessed by at least three annotators. Annotators verified label-image alignment, measurement units, bounding boxes, and report text consistency. Disagreements were resolved via majority voting to ensure annotation correctness and clinical plausibility.

Task-Specific Prompt Designing. To ensure consistent model behavior and fair comparability across tasks, we design structured prompts for each of the 50 application scenarios, consisting of three components: (1) a clinical role definition to set context and expertise, (2) a task-specific instruction aligned with standard sonography workflow, and (3) an output format specification, such as classification options, value ranges, or reference output examples. Detailed prompts are included in Appendix D. An ablation study on the impact of prompt design is presented in Section 5.1.

3.3 Statistics

U2-BENCH comprises 7,241 ultrasound studies spanning 8 benchmark tasks and 15 anatomical regions. Table 1 details the number of cases per task. Classification and detection constitute the largest shares, with 2,999 and 2,921 cases, respectively, while generation and regression tasks provide targeted evaluation of report synthesis and clinical value estimation.

Figure 2 summarizes the distribution across anatomical regions. Thyroid and breast ultrasound together account for more than one-third of all cases. This is because of their high clinical prevalence and broad diagnostic utility. Many anatomies support multiple tasks - for instance, fetal ultrasound is used for classification and regression - enabling multi-task evaluation within a unified anatomical context. This composition ensures broad coverage across modalities, tasks, and body regions, supporting robust and clinically grounded assessment of LVLM performance.

4 Experiment

4.1 Evaluation Settings

We conducted experiments on **U2-BENCH** with both open-source and closed-source LVLMs. Uniform prompts were applied across all models. The evaluation was executed on 32 NVIDIA A800 GPUs over a period of approximately two weeks, using the OpenCompass VLMEvalKit [23], with additional support from a unified framework [82].

Evaluated Models. We evaluated 20 LVLMs, spanning both open-source and closed-source systems, and including both general-purpose and medical-specialized variants.

- **Qwen2.5-VL Series [84]:** This includes *Qwen2.5-VL-3B-Instruct*, *Qwen2.5-VL-7B-Instruct*, *Qwen2.5-VL-32B-Instruct*, *Qwen2.5-VL-72B-Instruct*.
- **Medical-Specific Open-Source Models:** *MiniGPT-Med* [78], *MedDr* [26]
- **Other Open-Source Models:** *Phi-4-Multimodal-Instruct-5.6B* [1], *InternVL3-9B-Instruct* [90], *LLaVA-1.5-13B* [45], *Mistral-Small-3.1-24B-Instruct-2503* [35], *DeepSeek-VL2* [20]
- **Closed-Source Models:** *GPT-4o-Mini*, *GPT-4o-2024-08-06* [54], *Gemini-1.5-Pro (exp-02-05)* [4], *Gemini-2-Pro (exp-02-05)*, *Gemini-2.5-Pro-Preview (exp-02-05)* [4], *Claude-3-Sonnet (20250219)* [5], *Qwen-Max-2025-01-25* [8], *Doubao-1.5-Vision-Pro-32K-250115* [14], *Dolphin-V1* (Model developed by *Dolphin AI*)

4.2 Evaluation Protocol

We employed standard metrics aligned with clinical relevance and prior LVLM benchmarks. Classification tasks were evaluated with accuracy and F1 score. Detection tasks were converted to position classification tasks, and hence utilized accuracy as a metric to assess localization correctness. Regression tasks report Root Mean Squared Error (RMSE), Mean Absolute Error (MAE) and percentage within tolerance (%_tol). Generation tasks were assessed with BLEU-4 as percentage, ROUGE, and BERTScore [87] to capture both lexical and semantic similarity. All metrics were computed using ground-truth labels from the original dataset and standardized outputs with format specified by the prompts across models to ensure fair comparison.

U2-Score. We design a quantitative score to provide an overall evaluation metric for the ultrasound understanding capability of a model. The **U2-SCORE** is defined as a weighted combination of the metrics across all tasks. This can be formulated as:

$$\text{U2-Score} := \sum_{t=1}^N w_t d_t, \text{ where } w_t = \frac{n_t}{\sum_j n_j}, \text{ and } d_t \leq 1 \quad (1)$$

where N represents the number of tasks, w_t is the corresponding task weight, which is computed from the proportion of the sample number n_t of the t -th task. This can mitigate the imbalance issue of sample size in different tasks. Here, d_t denotes the value of the selected metric of the t -th task. Table 2 presents the values of w_t and the corresponding metrics being selected.

Table 2: **Task-specific evaluation metrics and weights.** The corresponding weight w_t and metric used for overall score aggregation for each task.

| t | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|--------|------|------|------|------|------|--------|--------|--------|
| | DD | VRA | LL | OD | KD | CVE | RG | CG |
| w_t | 0.2 | 0.2 | 0.07 | 0.27 | 0.07 | 0.07 | 0.08 | 0.04 |
| Metric | Acc. | Acc. | Acc. | Acc. | Acc. | 1-RMSE | BLEU-4 | BLEU-4 |

4.3 Evaluation Results

We present a comprehensive comparison of multimodal models on the **U2-BENCH** benchmark (Table 3), aiming to identify key performance trends across tasks and model types.

Closed-Source Models Lead. Closed-source models continue to dominate, with **Gemini-2.5-Pro-Preview** achieving the highest overall score of **0.2968**, narrowly surpassing the best-performing open-source model **DeepSeek-VL2** by a margin of just 0.0338. Other strong proprietary models such as **Dolphin-V1** and **Doubao-1.5-Vision-Pro** also perform competitively, highlighting that while open-source models are progressing rapidly, access to proprietary data and task-specific optimization still provides a measurable edge.

Task Difficulty Varies Significantly. Image classification tasks are generally easier, with **Doubao** reaching an accuracy of **0.558** on **DD**, and eight models—both open- and closed-source—exceeding 0.48. In contrast, spatial reasoning and text generation remain difficult: no model achieves accuracy above 0.16 on **KD**, and all models score below 8.0 BLEU on **RG**. Regression tasks (e.g., **CVE**) are also challenging, only the closed-source **Qwen-Max** reduces RMSE to 0.1248, while all open-source models remain above 0.1675.

Table 3: Results of different models on the **U2-BENCH**. We utilize **green** (1st), **blue** (2nd), and **yellow** (3rd) backgrounds to distinguish the top three results within different models. The “U2-Score” column represents the quantitative score defined in Section 4.2. To calculate the **U2-Score** for random guessing, the BLEU scores are taken to be zero.

| Models | DD | | VRA | | LL | | OD | | KD | | CVE | | RG | | | CG | | | U2-Score ↑ |
|-------------------------------------|--------|--------|--------|--------|--------|--------|--------|--------|--------|---------|--------|---------|---------|---------|----------|---------|---------|----------|------------|
| | Acc. ↑ | F1 ↑ | Acc. ↑ | F1 ↑ | Acc. ↑ | F1 ↑ | Acc. ↑ | F1 ↑ | Acc. ↑ | F1 ↑ | RMSE ↓ | MAE ↓ | %_tol ↑ | BLEU% ↑ | Rouge% ↑ | BERT% ↑ | BLEU% ↑ | Rouge% ↑ | |
| Random Guessing | 0.4143 | 0.4135 | 0.3195 | 0.3184 | 0.1118 | 0.0680 | 0.1120 | 0.5472 | 0.4352 | 18.776 | - | - | - | - | - | - | - | - | 0.2125 |
| Medical-Specific Open-Source Models | | | | | | | | | | | | | | | | | | | |
| MiniGPT-Med | 0.3468 | 0.2828 | 0.1800 | 0.1048 | 0.1728 | 0.1789 | 0.0840 | 0.3056 | 0.2600 | 33.2259 | 6.4700 | 20.1300 | 74.6900 | 30.2000 | 47.7500 | 80.5000 | 0.2375 | | 0.2375 |
| MedDr | 0.4508 | 0.3118 | 0.2071 | 0.1214 | 0.0720 | 0.0881 | 0.0900 | 0.2144 | 0.1786 | 38.2642 | 2.7998 | 13.5060 | 72.2050 | 33.4939 | 49.6236 | 81.2078 | 0.2373 | | 0.2373 |
| Open-Source Multimodal Models | | | | | | | | | | | | | | | | | | | |
| Qwen-2.5-VL-3B-Instruct | 0.4503 | 0.3591 | 0.2097 | 0.1492 | 0.0696 | 0.0649 | 0.0894 | 0.5008 | 0.4519 | 18.9055 | 3.5018 | 15.0327 | 72.8419 | 27.6748 | 44.7618 | 79.8849 | 0.2095 | | 0.2095 |
| Qwen-2.5-VL-7B-Instruct | 0.4821 | 0.3860 | 0.2181 | 0.1665 | 0.0750 | 0.0704 | 0.1000 | 0.4646 | 0.4337 | 19.7115 | 3.7100 | 15.5600 | 73.1500 | 29.4400 | 47.0000 | 81.1500 | 0.2235 | | 0.2235 |
| Qwen-2.5-VL-32B-Instruct | 0.4812 | 0.3860 | 0.2864 | 0.2071 | 0.1700 | 0.0755 | 0.0880 | 0.3414 | 0.3015 | 27.4038 | 1.9000 | 13.0100 | 68.1400 | 14.7700 | 38.6800 | 77.3900 | 0.2449 | | 0.2449 |
| Qwen-2.5-VL-72B-Instruct | 0.4895 | 0.4556 | 0.2559 | 0.1789 | 0.1150 | 0.0660 | 0.0860 | 0.3224 | 0.2733 | 37.9370 | 3.0900 | 15.0600 | 72.6600 | 28.1600 | 44.2800 | 80.9100 | 0.2421 | | 0.2421 |
| DeepSeek-VL2 | 0.4126 | 0.3190 | 0.2268 | 0.1111 | 0.2950 | 0.1682 | 0.1320 | 0.2956 | 0.2505 | 12.3355 | 7.4700 | 20.5400 | 75.3800 | 11.4200 | 34.8500 | 77.2400 | 0.2630 | | 0.2630 |
| InternVL3-9B-Instruct | 0.4447 | 0.3716 | 0.1926 | 0.1083 | 0.3000 | 0.1416 | 0.0940 | 0.2429 | 0.1733 | 50.8738 | 2.1600 | 14.7000 | 72.2100 | 21.5900 | 43.1300 | 80.9800 | 0.2566 | | 0.2566 |
| LLaVA-1.5-13B | 0.4321 | 0.3055 | 0.1731 | 0.0755 | 0.1700 | 0.1259 | 0.1100 | 0.2307 | 0.1976 | 24.7964 | 6.2400 | 18.5800 | 73.7900 | 10.8300 | 29.4000 | 75.5000 | 0.2378 | | 0.2378 |
| Phi-4-Multimodal-Instruct | 0.3686 | 0.1148 | 0.2452 | 0.0537 | 0.0350 | 0.0815 | 0.1600 | 0.2249 | 0.2006 | 16.1972 | 3.2700 | 16.5800 | 73.2700 | 3.8700 | 22.9800 | 73.0800 | 0.2168 | | 0.2168 |
| Mistral-Small-3.1-24B-Inst. | 0.4359 | 0.0936 | 0.1964 | 0.0664 | 0.1300 | 0.0910 | 0.1060 | 0.1675 | 0.1331 | 45.9459 | 1.8000 | 14.9000 | 71.7200 | 20.7700 | 42.1200 | 80.7400 | 0.2356 | | 0.2356 |
| Closed-Source Multimodal Models | | | | | | | | | | | | | | | | | | | |
| Doubao-1.5-Vision-Pro-32k | 0.5580 | 0.2597 | 0.2922 | 0.2147 | 0.1700 | 0.0729 | 0.1240 | 0.3664 | 0.3377 | 33.1731 | 0.7100 | 6.6450 | 72.4000 | 8.6400 | 33.3000 | 78.4200 | 0.2587 | | 0.2587 |
| GPT-4o-Mini | 0.4924 | 0.3784 | 0.1922 | 0.1272 | 0.1357 | 0.0846 | 0.0960 | 0.2267 | 0.1976 | 19.2308 | 4.9400 | 17.5200 | 74.1300 | 11.7300 | 36.2900 | 77.3300 | 0.2388 | | 0.2388 |
| GPT-4o | 0.4928 | 0.4132 | 0.1504 | 0.0974 | 0.1161 | 0.0850 | 0.0840 | 0.3712 | 0.3527 | 15.7895 | 2.6800 | 14.7700 | 73.3500 | 33.7700 | 49.9600 | 81.8800 | 0.2253 | | 0.2253 |
| Gemini-1.5-Pro | 0.3781 | 0.2247 | 0.0909 | 0.0476 | 0.2700 | 0.0661 | 0.0980 | 0.2772 | 0.2205 | 40.7051 | 0.5800 | 9.9400 | 70.5500 | 28.5800 | 45.9200 | 80.0200 | 0.1999 | | 0.1999 |
| Gemini-2.0-Pro-Exp | 0.4925 | 0.4194 | 0.1648 | 0.1323 | 0.1714 | 0.0945 | 0.0820 | 0.1945 | 0.1498 | 53.3333 | 0.2600 | 6.9200 | 40.2400 | 31.1800 | 48.6000 | 81.6000 | 0.2438 | | 0.2438 |
| Gemini-2.5-Pro-Preview | 0.4256 | 0.3112 | 0.2098 | 0.1493 | 0.2709 | 0.2714 | 0.2518 | 0.2937 | 0.2672 | 34.4970 | 5.5030 | 18.0180 | 74.4930 | 15.0110 | 38.0070 | 75.9890 | 0.2968 | | 0.2968 |
| Claude-3.7-Sonnet | 0.2121 | 0.0449 | 0.1453 | 0.0479 | 0.1356 | 0.0540 | 0.0760 | 0.1764 | 0.1500 | 36.0215 | 0.6900 | 12.2300 | 68.7400 | 1.2900 | 16.6600 | 71.6600 | 0.1596 | | 0.1596 |
| Qwen-Max | 0.4566 | 0.2676 | 0.1925 | 0.0871 | 0.1606 | 0.0761 | 0.0940 | 0.1248 | 0.0843 | 69.2308 | 3.5000 | 17.0200 | 73.9600 | 30.6700 | 49.0000 | 82.5500 | 0.2445 | | 0.2445 |
| Dolphin-V1 | 0.5107 | 0.4173 | 0.3406 | 0.2181 | 0.1950 | 0.0791 | 0.1500 | 0.1898 | 0.1463 | 56.2500 | 0.9300 | 11.5400 | 71.0600 | 27.2800 | 43.8600 | 80.0800 | 0.2841 | | 0.2841 |

Scaling Brings Diminishing Returns. Within the **Qwen-2.5-VL** family, scaling from 3B to 72B parameters yields consistent performance gains. While larger models achieve lower **CVE** RMSE, improvements in language generation and spatial reasoning tasks plateau, suggesting that excessive scaling may lead to overfitting on superficial visual patterns, ultimately harming clinical text generation capabilities.

Domain-Specific Models Excel in Reasoning. Medical-domain models such as **MedDr** demonstrate strong performance on clinical reasoning tasks (CVE RMSE = 0.214; CG BERT = 81.21), outperforming many general-purpose systems. However, they lag behind general multimodal models on visual classification (e.g., Qwen-72B DD F1 = 0.456). This highlights that domain-specialized models are better suited for semantic and quantitative tasks, while general models still excel at coarse-grained visual understanding. Combining both may offer a promising direction for improving overall performance.

5 Analysis

5.1 Causal Analysis of Prompt Designing

We investigate whether explicitly naming the anatomical region in the prompt has a *causal effect* on LVLm diagnostic accuracy in ultrasound.

Formal Setup. We examine the causal relationship between input image X , prompt P and model output A . We assume the prompt has two components P_1 (general task context) and P_2 (anatomy token), and construct the causal graph as shown in Figure 4. In this ablation study, we focus on whether the inclusion of anatomical information in the prompt causally affects the model output A . That is, whether P_2 has a causal effect on A , denoted by the question-marked arrow.

Anatomy Information Manipulation. Following the setup from [86], the effect of the confounders can be ignored and we utilize Pearl’s front-door adjustment. That is, we want to examine $\mathbb{P}(A|do(P_2))$ where P_2 gives the information about the anatomy. Then the causal effect of

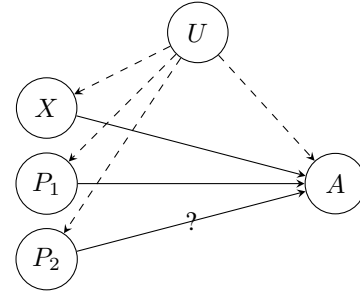


Figure 4: **Causal graph** of X , P_1 , and P_2 influence A , with a potential confounder U affecting all nodes. We examine the causal effect of $P_2 \rightarrow A$.

P_2 can be expressed as:

$$\mathbb{P}(A \mid do(P_2 = p_2^*)) = \sum_x \sum_{p_1} \mathbb{P}(A \mid x, p_1, p_2^*) \mathbb{P}(x \mid p_2^*) \mathbb{P}(p_1) \quad (2)$$

where $\mathbb{P}(p_1)$ is fixed by our prompt template, and $\mathbb{P}(x \mid p_2^*)$ is controlled to be the same in the experiment. Therefore, it suffices to investigate the distribution difference between $A_1 = \mathcal{M}(X, p_1, 1)$ and $A_0 = \mathcal{M}(X, p_1, 0)$, where \mathcal{M} denotes the model, $p_2^* = 1$ represents token present and $p_2^* = 0$ represents token absent. For more details about the causal analysis, please refer to Appendix E.

Experiment. We perform a McNemar’s test on 521 samples of breast and thyroid from the dataset on model Gemini-2.0-Pro-Exp, to investigate the distribution of A_1 and A_0 . We define two prompt variants for each image:

With anatomy ($P_2 = 1$): “You are a radiologist analysing a {anatomy} ultrasound image...”
No anatomy ($P_2 = 0$): “You are a radiologist analysing an ultrasound image...”

Each image x_i is evaluated under both conditions $\{(x_i, p_1, 1), (x_i, p_1, 0)\}$. Each prompt–image pair is forward-passed 5 times and majority-voted into a final prediction.

Results. Table 4 shows the paired contingency counts for the *no-anatomy* (ablation) versus *with-anatomy* prompts on $N = 521$ studies.

Table 4: **Effect of anatomy tokens in prompt design.** Paired comparison of model outcomes on 521 samples using prompts with vs. without anatomy-specific tokens. Each cell shows the number of samples falling into the respective outcome combination. Including anatomy information improves overall accuracy by 7.3 percentage points.

| With-anatomy prompt | No-anatomy prompt | |
|---------------------|---------------------------------------|------------------------------------|
| | Correct | Incorrect |
| Correct | 209 (<i>both correct</i>) | 64 (<i>only anatomy correct</i>) |
| Incorrect | 26 (<i>only no-anatomy correct</i>) | 222 (<i>both incorrect</i>) |

* Net effect: Accuracy with anatomy = 52.4 %, without = 45.1 %; gain = +7.3 pp (McNemar $\chi^2 = 16.0$, $p = 6.2 \times 10^{-5}$).

The McNemar’s exact test yields a χ_1^2 value of 16.04, which corresponds to a p-value of 6.2×10^{-5} and suggests a significant statistical difference between the two conditions. Therefore, we conclude that there is strong evidence that including anatomy information in the prompt can improve the performance of the model.

5.2 Instruction Following Analysis

Table 5 shows that contemporary models are already highly adept at parsing prompts and adhering to output specifications: six of the seventeen systems achieve a perfect score on the DD benchmark. The remaining models lag only slightly behind. The medical-oriented MiniGPT-Med [3] and MedDr [26] deliver middling results, while Qwen-3B and Qwen-72B [9] close the gap rapidly as their parameter counts increase. Claude-3.7 [6] score of 0.942 is largely attributable to occasional formatting omissions. For every non-perfect model, the deviation from the maximum is under six percentage points, and no systematic failures are observed.

Table 5: Instruction following comparison across different models.

| Task | Models | | | | | | | | | | | | | | | | |
|------|-------------|-------|---------|---------|----------|----------|------------|----------|----------|-------|-------|---------|------------|--------|---------|---------|------------|
| | MiniGPT-Med | MedDr | Qwen-3B | Qwen-7B | Qwen-32B | Qwen-72B | Dolphin-V1 | DeepSeek | InternVL | LLaVA | Phi-4 | Mistral | Doubao-1.5 | GPT-4o | Gem-2.0 | Gem-2.5 | Claude-3.7 |
| DD | 0.952 | 0.961 | 0.968 | 0.983 | 0.996 | 1.000 | 1.000 | 1.000 | 0.993 | 0.987 | 0.998 | 0.999 | 1.000 | 1.000 | 0.997 | 1.000 | 0.942 |

6 Conclusion

Ultrasound is essential to global healthcare but remains difficult to interpret. We present **U2-BENCH**, the first benchmark for evaluating LVLMS on ultrasound understanding. It includes 7,241 cases across 15 anatomical regions and defines 8 clinical tasks for 50 application scenarios. Evaluating 20 LVLMS, we find their strong performance in classification but persistent challenges in spatial reasoning and clinical text generation, suggesting a future direction for improving LVLMS on ultrasound interpretation.

References

- [1] Abdin, M.I., Aneja, J., Behl, H.S., Bubeck, S., Eldan, R., Gunasekar, S., Harrison, M., Hewett, R.J., Javaheripi, M., Kauffmann, P., Lee, J.R., Lee, Y.T., Li, Y., Liu, W., Mendes, C.C.T., Nguyen, A., Price, E., de Rosa, G., Saarikivi, O., Salim, A., Shah, S., Wang, X., Ward, R., Wu, Y., Yu, D., Zhang, C., Zhang, Y., 2024. Phi-4 technical report. CoRR abs/2412.08905. URL: <https://doi.org/10.48550/arXiv.2412.08905>, doi:10.48550/ARXIV.2412.08905, arXiv:2412.08905.
- [2] Al-Dhabyani, W., Gomaa, M., Khaled, H., Fahmy, A., 2020. Dataset of breast ultrasound images. Data in Brief 28, 104863. URL: <https://www.sciencedirect.com/science/article/pii/S2352340919312181>, doi:<https://doi.org/10.1016/j.dib.2019.104863>.
- [3] Alkhalidi, A., Alnajim, R., Alabdullatef, L., Alyahya, R., Chen, J., Zhu, D., Alsinan, A., Elhoseiny, M., 2024. Minigpt-med: Large language model as a general interface for radiology diagnosis. arXiv preprint arXiv:2407.04106 .
- [4] Anil, R., Borgeaud, S., Wu, Y., Alayrac, J., Yu, J., Soricut, R., Schalkwyk, J., Dai, A.M., Hauth, A., Millican, K., Silver, D., Petrov, S., Johnson, M., Antonoglou, I., Schrittwieser, J., Glaese, A., Chen, J., Pitler, E., Lillicrap, T.P., Lazaridou, A., Firat, O., Molloy, J., Isard, M., Barham, P.R., Hennigan, T., Lee, B., Viola, F., Reynolds, M., Xu, Y., Doherty, R., Collins, E., Meyer, C., Rutherford, E., Moreira, E., Ayoub, K., Goel, M., Tucker, G., Piqueras, E., Krikun, M., Barr, I., Savinov, N., Danihelka, I., Roelofs, B., White, A., Andreassen, A., von Glehn, T., Yagati, L., Kazemi, M., Gonzalez, L., Khalman, M., Sygnowski, J., et al., 2023. Gemini: A family of highly capable multimodal models. CoRR abs/2312.11805.
- [5] Anthropic, 2024. Claude 3.5 sonnet. Anthropic News URL: <https://www.anthropic.com/news/claude-3-5-sonnet>.
- [6] Anthropic, 2025. Claude 3.7 sonnet. Anthropic News URL: <https://www.anthropic.com/news/claude-3-7-sonnet>.
- [7] Bai, J., 2024. Psfhs. URL: <https://doi.org/10.5281/zenodo.10969427>, doi:10.5281/zenodo.10969427.
- [8] Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., Huang, F., Hui, B., Ji, L., Li, M., Lin, J., Lin, R., Liu, D., Liu, G., Lu, C., Lu, K., Ma, J., Men, R., Ren, X., Ren, X., Tan, C., Tan, S., Tu, J., Wang, P., Wang, S., Wang, W., Wu, S., Xu, B., Xu, J., Yang, A., Yang, H., Yang, J., Yang, S., Yao, Y., Yu, B., Yuan, H., Yuan, Z., Zhang, J., Zhang, X., Zhang, Y., Zhang, Z., Zhou, C., Zhou, J., Zhou, X., Zhu, T., 2023a. Qwen technical report. arXiv preprint arXiv:2309.16609 .
- [9] Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., Zhou, J., 2023b. Qwen-vl: A frontier large vision-language model with versatile abilities. CoRR abs/2308.12966.
- [10] Ben Abacha, A., Hasan, S.A., Datla, V.V., Liu, J., Demner-Fushman, D., Müller, H., 2019. Vqa-med: Overview of the medical visual question answering task at imageclef 2019, in: Working Notes of CLEF 2019, CEUR-WS.org, Lugano, Switzerland. URL: https://ceur-ws.org/Vol-2380/paper_272.pdf.
- [11] Bi, Y., Jiang, Z., Clarenbach, R., Ghotbi, R., Karlas, A., Navab, N., 2024. Mi-segnet: Mutual information-based us segmentation for unseen domain generalization. URL: <https://arxiv.org/abs/2303.12649>, arXiv:2303.12649.
- [12] Burgos-Artizzu, X.P., Coronado-Gutierrez, D., Valenzuela-Alcaraz, B., Bonet-Carne, E., Eixarch, E., Crispi, F., Gratacós, E., 2020. Fetal_planes_db: Common maternal-fetal ultrasound images. URL: <https://doi.org/10.5281/zenodo.3904280>, doi:10.5281/zenodo.3904280.
- [13] Byra, M., Styczynski, G., Szmigielski, C., Kalinowski, P., Michalowski, L., Paluszkiwicz, R., Ziarkiewicz-Wroblewska, B., Zieniewicz, K., Sobieraj, P., Nowicki, A., 2018. Dataset of b-mode fatty liver ultrasound images. URL: <https://doi.org/10.5281/zenodo.1009146>, doi:10.5281/zenodo.1009146.

- [14] ByteDance, 2024. Doubao model series. <https://www.doubao.com/>. Accessed: 2025-05-14.
- [15] Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H.P.d.O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., et al., 2021. Evaluating large language models trained on code. arXiv preprint arXiv:2107.03374 .
- [16] Chen, P., Ye, J., Wang, G., Li, Y., Deng, Z., Li, W., Li, T., Duan, H., Huang, Z., Su, Y., Wang, B., Zhang, S., Fu, B., Cai, J., Zhuang, B., Seibel, E.J., He, J., Qiao, Y., 2024a. Gmai-mmbench: A comprehensive multimodal evaluation benchmark towards general medical ai. URL: <https://arxiv.org/abs/2408.03361>, arXiv:2408.03361.
- [17] Chen, P., Ye, J., Wang, G., Li, Y., Deng, Z., Li, W., Li, T., Duan, H., Huang, Z., Su, Y., Wang, B., Zhang, S., Fu, B., Cai, J., Zhuang, B., Seibel, E.J., Qiao, Y., He, J., 2024b. Gmai-mmbench: A comprehensive multimodal evaluation benchmark towards general medical ai, in: Globerson, A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J., Zhang, C. (Eds.), Advances in Neural Information Processing Systems, Curran Associates, Inc.. pp. 94327–94427. URL: https://proceedings.neurips.cc/paper_files/paper/2024/file/ab7e02fd60e47e2a379d567f6b54f04e-Paper-Datasets_and_Benchmarks_Track.pdf.
- [18] Chen, Z., Wu, J., Zhou, J., Wen, B., Bi, G., Jiang, G., Cao, Y., Hu, M., Lai, Y., Xiong, Z., Huang, M., 2024c. Tombench: Benchmarking theory of mind in large language models, in: Ku, L., Martins, A., Srikumar, V. (Eds.), Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024, Association for Computational Linguistics. pp. 15959–15983.
- [19] Da Correggio, K.S., Noya Galluzzo, R., Santos, L.O., Soares Muylaert Barroso, F., Zimmermann Loureiro Chaves, T., Sherley Casimiro Onofre, A., von Wangenheim, A., 2023. Fetal abdominal structures segmentation dataset using ultrasonic images. URL: <https://doi.org/10.17632/4gcpm9dsc3.1>, doi:10.17632/4gcpm9dsc3.1.
- [20] DeepSeek-AI, 2024. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. arXiv:2405.04434.
- [21] DeepSeek-AI, Liu, A., Feng, B., Xue, B., Wang, B., Wu, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, C., Dai, D., Guo, D., Yang, D., Chen, D., Ji, D., Li, E., Lin, F., Dai, F., Luo, F., Hao, G., Chen, G., Li, G., Zhang, H., Bao, H., Xu, H., Wang, H., Zhang, H., Ding, H., Xin, H., Gao, H., Li, H., Qu, H., Cai, J.L., Liang, J., Guo, J., Ni, J., Li, J., Wang, J., Chen, J., Chen, J., Yuan, J., Qiu, J., Li, J., Song, J., Dong, K., Hu, K., Gao, K., Guan, K., Huang, K., Yu, K., Wang, L., Zhang, L., Xu, L., Xia, L., Zhao, L., Wang, L., Zhang, L., Li, M., Wang, M., Zhang, M., Zhang, M., Tang, M., Li, M., Tian, N., Huang, P., Wang, P., Zhang, P., Wang, Q., Zhu, Q., Chen, Q., Du, Q., Chen, R.J., Jin, R.L., Ge, R., Zhang, R., Pan, R., Wang, R., Xu, R., Zhang, R., Chen, R., Li, S.S., Lu, S., Zhou, S., Chen, S., Wu, S., Ye, S., Ye, S., Ma, S., Wang, S., Zhou, S., Yu, S., Zhou, S., Pan, S., Wang, T., Yun, T., Pei, T., Sun, T., Xiao, W.L., Zeng, W., 2024. Deepseek-v3 technical report. CoRR abs/2412.19437. URL: <https://doi.org/10.48550/arXiv.2412.19437>, doi:10.48550/ARXIV.2412.19437, arXiv:2412.19437.
- [22] Divekar, A., Sonawane, A., 2024. Leveraging ai for automatic classification of pcos using ultrasound imaging. URL: <https://arxiv.org/abs/2501.01984>, arXiv:2501.01984.
- [23] Duan, H., Yang, J., Qiao, Y., Fang, X., Chen, L., Liu, Y., Dong, X., Zang, Y., Zhang, P., Wang, J., et al., 2024. Vlmevalkit: An open-source toolkit for evaluating large multi-modality models, in: Proceedings of the 32nd ACM International Conference on Multimedia, pp. 11198–11201.
- [24] FALLMUD, 2025. FALLMUD: Fascicle lower leg muscle ultrasound dataset. <https://kalisteo.cea.fr/index.php/fallmud/>. Dataset composed of 812 lower leg muscle ultrasound images with segmentation masks, used for muscle structure analysis and injury prevention.
- [25] Guo, X., Men, Q., Noble, J.A., 2024. MMSummary: Multimodal Summary Generation for Fetal Ultrasound Video. URL: <http://arxiv.org/abs/2408.03761>, doi:10.48550/arXiv.2408.03761. arXiv:2408.03761 [cs].

- [26] He, S., Nie, Y., Wang, H., Yang, S., Wang, Y., Cai, Z., Chen, Z., Xu, Y., Luo, L., Xiang, H., Lin, X., Wu, M., Peng, Y., Shih, G., Xu, Z., Wu, X., Wang, Q., Chan, R.C.K., Vardhanabhuti, V., Chu, W.C.W., Zheng, Y., Rajpurkar, P., Zhang, K., Chen, H., 2024. Gsco: Towards generalizable ai in medicine via generalist-specialist collaboration. URL: <https://arxiv.org/abs/2404.15127>, arXiv:2404.15127.
- [27] He, X., Zhang, Y., Mou, L., Xing, E., Xie, P., 2020. Pathvqa: 30000+ questions for medical visual question answering. arXiv preprint arXiv:2003.10286 .
- [28] van den Heuvel, T.L.A., de Bruijn, D., de Korte, C.L., van Ginneken, B., 2018. Automated measurement of fetal head circumference using 2d ultrasound images. URL: <https://doi.org/10.5281/zenodo.1327317>, doi:10.5281/zenodo.1327317.
- [29] Hewson, D.W., Bedforth, N.M., 2023. Closing the gap: artificial intelligence applied to ultrasound-guided regional anaesthesia. *British Journal of Anaesthesia* 130, 245–247. URL: <https://www.sciencedirect.com/science/article/pii/S0007091222006924>, doi:<https://doi.org/10.1016/j.bja.2022.12.005>.
- [30] Huang, W., Abbeel, P., Pathak, D., Mordatch, I., 2022. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. arXiv preprint arXiv: Arxiv-2201.07207 .
- [31] Huang, X., Shen, L., Liu, J., Shang, F., Li, H., Huang, H., Yang, Y., 2025. Towards a multimodal large language model with pixel-level insight for biomedicine. *Proceedings of the AAAI Conference on Artificial Intelligence* 39, 3779–3787. URL: <https://ojs.aaai.org/index.php/AAAI/article/view/32394>, doi:10.1609/aaai.v39i4.32394.
- [32] Huang, Z., Bianchi, F., Yuksekgonul, M., Montine, T.J., Zou, J., 2023. A visual–language foundation model for pathology image analysis using medical twitter. *Nature Medicine* 29, 2307–2316. URL: <https://doi.org/10.1038/s41591-023-02504-3>, doi:10.1038/s41591-023-02504-3.
- [33] Indelman, H.C., Dahan, E., Perez-Agosto, A.M., Shiran, C., Shaked, D., Daniel, N., 2024. Semantic Segmentation Refiner for Ultrasound Applications with Zero-Shot Foundation Models. URL: <http://arxiv.org/abs/2404.16325>, doi:10.48550/arXiv.2404.16325. arXiv:2404.16325 [cs] version: 1.
- [34] Ji, Y., Bai, H., GE, C., Yang, J., Zhu, Y., Zhang, R., Li, Z., Zhanng, L., Ma, W., Wan, X., Luo, P., 2022. Amos: A large-scale abdominal multi-organ benchmark for versatile medical image segmentation, in: Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., Oh, A. (Eds.), *Advances in Neural Information Processing Systems*, Curran Associates, Inc., pp. 36722–36732. URL: https://proceedings.neurips.cc/paper_files/paper/2022/file/ee604e1bedbd069d9fc9328b7b9584be-Paper-Datasets_and_Benchmarks.pdf.
- [35] Jiang, A.Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D.S., Casas, D.d.l., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., et al., 2023. Mistral 7b. arXiv preprint arXiv:2310.06825 .
- [36] Krönke, M., Eilers, C., Dimova, D., Köhler, M., Buschner, G., Schweiger, L., Konstantinidou, L., Makowski, M., Nagarajah, J., Navab, N., Weber, W., Wendler, T., 2022. Tracked 3d ultrasound and deep neural network-based thyroid segmentation reduce interobserver variability in thyroid volumetry. *PLOS ONE* 17, 1–15. URL: <https://doi.org/10.1371/journal.pone.0268550>, doi:10.1371/journal.pone.0268550.
- [37] Lau, J.J., Gayen, S., Ben Abacha, A., Demner-Fushman, D., 2018. A dataset of clinically generated visual questions and answers about radiology images. *Scientific Data* 5, 180251. URL: <https://doi.org/10.1038/sdata.2018.251>, doi:10.1038/sdata.2018.251.
- [38] Laverde Saad, A., Jfri, A., García, R., Salguero, I., Martínez, C., Cembrero, H., Roustán, G., Alfageme, F., 2021. Discriminative deep learning based benignity/malignancy diagnosis of dermatologic ultrasound skin lesions with pretrained artificial intelligence architecture. *Skin Research and Technology* 28. doi:10.1111/srt.13086.

- [39] Leclerc, S., Smistad, E., Pedrosa, J., Østvik, A., Cervenansky, F., Espinosa, F., Espeland, T., Berg, E.A.R., Jodoin, P.M., Grenier, T., Lartizien, C., D’hooge, J., Lovstakken, L., Bernard, O., 2019. Deep learning for segmentation using an open large-scale dataset in 2d echocardiography. *IEEE Transactions on Medical Imaging* 38, 2198–2210. doi:10.1109/TMI.2019.2900516.
- [40] Leenings, R., Konowski, M., Winter, N.R., Ernsting, J., Fisch, L., Barkhau, C., Dannlowski, U., Lügering, A., Jiang, X., Hahn, T., 2025. C-trus: A novel dataset and initial benchmark for colon wall segmentation in transabdominal ultrasound, in: Gomez, A., Khanal, B., King, A., Namburete, A. (Eds.), *Simplifying Medical Ultrasound*, Springer Nature Switzerland, Cham. pp. 101–111.
- [41] Li, B., Wang, R., Wang, G., Ge, Y., Ge, Y., Shan, Y., 2023a. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*.
- [42] Li, J., Su, T., Zhao, B., Lv, F., Wang, Q., Navab, N., Hu, Y., Jiang, Z., 2024. Ultrasound report generation with cross-modality feature alignment via unsupervised guidance. URL: <https://arxiv.org/abs/2406.00644>, arXiv:2406.00644.
- [43] Li, J., Zhang, P., Wang, T., Wang, K., Sheng, B., 2023b. Lepset. URL: <https://doi.org/10.5281/zenodo.8041285>, doi:10.5281/zenodo.8041285.
- [44] Lin, Z., Lin, J., Zhu, L., Fu, H., Qin, J., Wang, L., 2022. A new dataset and a baseline model for breast lesion detection in ultrasound videos, in: Wang, L., Dou, Q., Fletcher, P.T., Speidel, S., Li, S. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2022*, Springer Nature Switzerland, Cham. pp. 614–623.
- [45] Liu, H., Li, C., Wu, Q., Lee, Y.J., 2023a. Visual instruction tuning.
- [46] Liu, X., Yu, H., Zhang, H., Xu, Y., Lei, X., Lai, H., Gu, Y., Ding, H., Men, K., Yang, K., et al., 2023b. AgentBench: Evaluating LLMs as agents. *arXiv preprint arXiv:2308.03688*.
- [47] Liu, Y., Duan, H., Zhang, Y., Li, B., Zhang, S., Zhao, W., Yuan, Y., Wang, J., He, C., Liu, Z., Chen, K., Lin, D., 2023c. Mmbench: Is your multi-modal model an all-around player? *arXiv:2307.06281*.
- [48] Lu, Y., Zhou, M., Zhi, D., Zhou, M., Jiang, X., Qiu, R., Ou, Z., Wang, H., Qiu, D., Zhong, M., Lu, X., Chen, G., Bai, J., 2022. The jnu-ifm dataset for segmenting pubic symphysis-fetal head. *Data in Brief* 41, 107904. URL: <https://www.sciencedirect.com/science/article/pii/S2352340922001160>, doi:<https://doi.org/10.1016/j.dib.2022.107904>.
- [49] Maroua, A., 2020. Algerian ultrasound images thyroid dataset (auitd). <https://www.kaggle.com/datasets/azouzmaroua/algeria-ultrasound-images-thyroid-dataset-auitd>. Kaggle dataset containing 3-class thyroid ultrasound images (Benign, Malignant, Normal) collected from hospitals in Setif, Algeria. Accessed May 2025.
- [50] Meiburger, K.M., Zahnd, G., Faita, F., Loizou, C., Carvalho, C., Steinman, D., Gibello, L., Bruno, R.M., Marzola, F., Clarenbach, R., Francesconi, M., Nicolaides, A., Campilho, A., Ghotbi, R., Kyriacou, E., Navab, N., Griffin, M., Panayiotou, A., Gherardini, R., Varetto, G., Bianchini, E., Pattichis, C., Ghiadoni, L., Rouco, J., Molinari, F., 2021. DATASET for "Carotid Ultrasound Boundary Study (CUBS): an open multi-center analysis of computerized intima-media thickness measurement systems and their clinical impact". URL: <https://doi.org/10.17632/fpv535fss7.1>, doi:10.17632/fpv535fss7.1.
- [51] Mo, Y., Han, C., Liu, Y., Liu, M., Shi, Z., Lin, J., Zhao, B., Huang, C., Qiu, B., Cui, Y., Wu, L., Pan, X., Xu, Z., Huang, X., Liu, Z., Wang, Y., Liang, C., 2022. Hover-trans: Anatomy-aware hover-transformer for roi-free breast cancer diagnosis in ultrasound images. URL: <https://arxiv.org/abs/2205.08390>, arXiv:2205.08390.
- [52] NeuronXJTU, palkia1998, 2023. Kfgnet: Source code for video classification using ultrasonic data. <https://github.com/NeuronXJTU/KFGNet>. GitHub repository, accessed May 2025. Includes Baidu Netdisk link to ultrasonic data.

- [53] Novin, S., Alvarez, C., Renner, J.B., Golightly, Y.M., Nelson, A.E., 2023. Features of knee and multijoint osteoarthritis by sex and race and ethnicity: A preliminary analysis in the johnston county health study. *Journal of Rheumatology* , jrheum.2023-0479doi:10.3899/jrheum.2023-0479. epub ahead of print, published September 15, 2023.
- [54] OpenAI, 2023. Gpt-4 technical report. PREPRINT .
- [55] Pahuni Choudhary, C., 2023. Carotid artery ultrasound and color doppler: Multimodal carotid artery imaging with comprehensive patient health records. <https://www.kaggle.com/datasets/pahunichoudhary/carotid-artery-ultrasound-and-color-doppler>. Kaggle dataset. Includes ultrasound and color Doppler images from 18 patients with structured health records. Licensed under Apache 2.0. Accessed May 2025.
- [56] Pawłowska, A., Ćwierz-Pieńkowska, A., Domalik, A., Jaguś, D., Kasprzak, P., Matkowski, R., Fura, Ł., Nowicki, A., Żołek, N., 2024. Curated benchmark dataset for ultrasound based breast lesion analysis. *Scientific Data* 11, 148. URL: <https://doi.org/10.1038/s41597-024-02984-z>, doi:10.1038/s41597-024-02984-z.
- [57] Pedraza, L., Vargas, C., Narváez, F., Durán, O., Muñoz, E., Romero, E., 2015. An open access thyroid ultrasound image database, in: 10th International symposium on medical information processing and analysis, SPIE. pp. 188–193.
- [58] Prabakaran, B., Hamelmann, P., Ostrowski, E., Shafique, M., 2023. Fpus23: An ultrasound fetus phantom dataset with deep neural network evaluations for fetus orientations, fetal planes, and anatomical features. *IEEE Access* PP, 1–1. doi:10.1109/ACCESS.2023.3284315.
- [59] Ravishankar, H., Patil, R., Melapudi, V., Suthar, H., Anzengruber, S., Bhatia, P., Taha, K.H., Annangi, P., 2023. SonoSAMTrack – Segment and Track Anything on Ultrasound Images. URL: <http://arxiv.org/abs/2310.16872>, doi:10.48550/arXiv.2310.16872. arXiv:2310.16872 [eess].
- [60] Sairam, V.A., 2020. Ultrasound breast images for breast cancer. <https://www.kaggle.com/datasets/aryashah2k/ultrasound-breast-images-for-breast-cancer>. Kaggle dataset. CC0: Public Domain license. Accessed May 2025.
- [61] Sappia, M.S., 2024. Acouslic-ai : Abdominal circumference operator- agnostic ultrasound measurement in low-income countries using artificial intelligence. URL: <https://doi.org/10.5281/zenodo.12697994>, doi:10.5281/zenodo.12697994.
- [62] Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., et al., 2022. LAION-5B: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems* 35, 25278–25294.
- [63] Sendra-Balcells, C., Campello, V.M., Torrents-Barrena, J., Ahmed, Y.A., Elattar, M., Ohene-Botwe, B., Nyangulu, P., Stones, W., Ammar, M., Benamer, L.N., Kitembo, H.N., Sereke, S.G., Wanyonyi, S.Z., Temmerman, M., Gratacós, E., Bonet, E., Eixarch, E., Mikolaj, K., Tolsgaard, M.G., Lekadir, K., 2023. Generalisability of fetal ultrasound deep learning models to low-resource imaging settings in five african countries. *Scientific Reports* 13, 2728. URL: <https://doi.org/10.1038/s41598-023-29490-3>, doi:10.1038/s41598-023-29490-3.
- [64] Shao, W., Brisbane, W., 2024. Micro-ultrasound prostate segmentation dataset. URL: <https://doi.org/10.5281/zenodo.10475293>, doi:10.5281/zenodo.10475293.
- [65] Sharma, H., Drukker, L., Papageorgiou, A.T., Noble, J.A., 2021. Machine learning-based analysis of operator pupillary response to assess cognitive workload in clinical ultrasound imaging. *Computers in Biology and Medicine* 135, 104589. URL: <https://www.sciencedirect.com/science/article/pii/S0010482521003838>, doi:https://doi.org/10.1016/j.combiomed.2021.104589.
- [66] Sharma, P., Ding, N., Goodman, S., Soricut, R., 2018. Conceptual Captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning, in: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2556–2565.

- [67] Shun-Shin, M., 2023. The imperial ai echocardiography dataset. <https://unityimaging.net>. An open-access dataset of over 7500 expert-labeled echocardiographic images for AI development in cardiology. IRAS: 279328, REC: 20/SC/0386. Accessed May 2025.
- [68] Singla, R., Ringstrom, C., Hu, G., Lessoway, V., Reid, J., Nguan, C., Rohling, R., 2023. The open kidney ultrasound data set, in: Kainz, B., Noble, A., Schnabel, J., Khanal, B., Müller, J.P., Day, T. (Eds.), *Simplifying Medical Ultrasound*, Springer Nature Switzerland, Cham. pp. 155–164.
- [69] Sivasubramaniam, S., Osei-Akoto, C., Zhang, Y., Stockinger, K., Fürst, J., 2024. Sm3-text-to-query: Synthetic multi-model medical text-to-query benchmark, in: Globerson, A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J., Zhang, C. (Eds.), *Advances in Neural Information Processing Systems*, Curran Associates, Inc.. pp. 88627–88663. URL: https://proceedings.neurips.cc/paper_files/paper/2024/file/a182a8e6ebc91728b6e6b6382c9f7b1e-Paper-Datasets_and_Benchmarks_Track.pdf.
- [70] Sun, L., Han, Y., Zhao, Z., Ma, D., Shen, Z., Chen, B., Chen, L., Yu, K., 2023. SciEval: A multi-level large language model evaluation benchmark for scientific research. arXiv preprint arXiv:2308.13149 .
- [71] Team, G., 2024. Advancing multimodal medical capabilities of gemini. CoRR abs/2405.03162.
- [72] Tyagi, A., Tyagi, A., Kaur, M., Aggarwal, R., Soni, K.D., Sivaswamy, J., Trikha, A., 2024. Nerve block target localization and needle guidance for autonomous robotic ultrasound guided regional anesthesia. URL: <https://arxiv.org/abs/2308.03717>, arXiv:2308.03717.
- [73] Vitale, S., Orlando, J.I., Iarussi, E., Larrabide, I., 2020. Improving realism in patient-specific abdominal ultrasound simulation using cyclegans. *International Journal of Computer Assisted Radiology and Surgery* 15, 183–192. doi:10.1007/s11548-019-02046-5. epub 2019 Aug 7.
- [74] Wang, Yipei and Yang, Q., Drukker, L., Papageorgiou, A., Hu, Y., Noble, J.A., 2022. Task model-specific operator skill assessment in routine fetal ultrasound scanning. *International Journal of Computer Assisted Radiology and Surgery* 17, 1437–1444. URL: <https://doi.org/10.1007/s11548-022-02642-y>, doi:10.1007/s11548-022-02642-y.
- [75] Wang, X., Hu, Z., Lu, P., Zhu, Y., Zhang, J., Subramaniam, S., Loomba, A.R., Zhang, S., Sun, Y., Wang, W., 2023. SciBench: Evaluating college-level scientific problem-solving abilities of large language models. arXiv preprint arXiv:2307.10635 .
- [76] Wiedemann, N., Boer, D.d.K.d., Richter, M., van de Weijer, S., Buhre, C., Eggert, F.A.M., Aarnoudse, S., Grevendonk, L., Röber, S., Remie, C.M., Buhre, W., Henry, R., Born, J., 2025. Covid-blues - a prospective study on the value of ai in lung ultrasound analysis. *IEEE Journal of Biomedical and Health Informatics* , 1–12doi:10.1109/JBHI.2025.3543686.
- [77] Wisesty, U.N., Thufailah, I.F., Dewi, R.M., Adiwijaya, Jondri, 2018. Study of segmentation technique and stereology to detect pco follicles on usg images. *Journal of Computer Science* 14, 351–359. URL: <https://thescipub.com/abstract/jcssp.2018.351.359>, doi:10.3844/jcssp.2018.351.359.
- [78] Wu, C., Zhang, X., Zhang, Y., Wang, Y., Xie, W., 2023a. Towards Generalist Foundation Model for Radiology by Leveraging Web-scale 2D&3D Medical Data. URL: <http://arxiv.org/abs/2308.02463>, doi:10.48550/arXiv.2308.02463. arXiv:2308.02463 [cs].
- [79] Wu, C., Zhang, X., Zhang, Y., Wang, Y., Xie, W., 2023b. Towards generalist foundation model for radiology by leveraging web-scale 2d&3d medical data. URL: <https://arxiv.org/abs/2308.02463>, arXiv:2308.02463.
- [80] Xia, P., Chen, Z., Tian, J., Gong, Y., Hou, R., Xu, Y., Wu, Z., Fan, Z., Zhou, Y., Zhu, K., Zheng, W., Wang, Z., Wang, X., Zhang, X., Bansal, C., Niethammer, M., Huang, J., Zhu, H., Li, Y., Sun, J., Ge, Z., Li, G., Zou, J., Yao, H., 2024. Cares: A comprehensive benchmark of trustworthiness in medical vision language models, in: Globerson,

- A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J., Zhang, C. (Eds.), *Advances in Neural Information Processing Systems*, Curran Associates, Inc.. pp. 140334–140365. URL: https://proceedings.neurips.cc/paper_files/paper/2024/file/fde7f40f8ced5735006810534dc66b33-Paper-Datasets_and_Benchmarks_Track.pdf.
- [81] Xiao, X., Zhang, J., Shao, Y., Liu, J., Shi, K., He, C., Kong, D., 2025. Deep learning-based medical ultrasound image and video segmentation methods: Overview, frontiers, and challenges. *Sensors* 25. URL: <https://www.mdpi.com/1424-8220/25/8/2361>, doi:10.3390/s25082361.
 - [82] XiaohuMini, 2025. Xiaohumini. XiaohuMini News URL: <https://xiaohumini.site>.
 - [83] Xu, Q., Hong, F., Li, B., Hu, C., Chen, Z., Zhang, J., 2023. On the tool manipulation capability of open-source large language models. *arXiv preprint arXiv: 2305.16504*.
 - [84] Yang, A., Yang, B., Hui, B., Zheng, B., Yu, B., Zhou, C., Li, C., Li, C., Liu, D., Huang, F., et al., 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.
 - [85] Ying, K., Meng, F., Wang, J., Li, Z., Lin, H., Yang, Y., Zhang, H., Zhang, W., Lin, Y., Liu, S., Lei, J., Lu, Q., Chen, R., Xu, P., Zhang, R., Zhang, H., Gao, P., Wang, Y., Qiao, Y., Luo, P., Zhang, K., Shao, W., 2024. Mmt-bench: A comprehensive multimodal benchmark for evaluating large vision-language models towards multitask agi. *arXiv:2404.16006*.
 - [86] Zhang, C., Zhang, L., Wu, J., He, Y., Zhou, D., 2024. Causal Prompting: Debiasing Large Language Model Prompting based on Front-Door Adjustment. URL: <http://arxiv.org/abs/2403.02738>, doi:10.48550/arXiv.2403.02738. *arXiv:2403.02738 [cs]*.
 - [87] Zhang*, T., Kishore*, V., Wu*, F., Weinberger, K.Q., Artzi, Y., 2020. Bertscore: Evaluating text generation with bert, in: *International Conference on Learning Representations*. URL: <https://openreview.net/forum?id=SkeHuCVFDr>.
 - [88] Zhong, W., Cui, R., Guo, Y., Liang, Y., Lu, S., Wang, Y., Saied, A., Chen, W., Duan, N., 2023. Agieval: A human-centric benchmark for evaluating foundation models. *arXiv:2304.06364*.
 - [89] Zhu, D., Chen, J., Shen, X., Li, X., Elhoseiny, M., 2023. MiniGPT-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.
 - [90] Zhu, J., Wang, W., Chen, Z., Liu, Z., Ye, S., Gu, L., Tian, H., Duan, Y., Su, W., Shao, J., Gao, Z., Cui, E., Wang, X., Cao, Y., Liu, Y., Wei, X., Zhang, H., Wang, H., Xu, W., Li, H., Wang, J., Deng, N., Li, S., He, Y., Jiang, T., Luo, J., Wang, Y., He, C., Shi, B., Zhang, X., Shao, W., He, J., Xiong, Y., Qu, W., Sun, P., Jiao, P., Lv, H., Wu, L., Zhang, K., Deng, H., Ge, J., Chen, K., Wang, L., Dou, M., Lu, L., Zhu, X., Lu, T., Lin, D., Qiao, Y., Dai, J., Wang, W., 2025. Internv13: Exploring advanced training and test-time recipes for open-source multimodal models. URL: <https://arxiv.org/abs/2504.10479>, *arXiv:2504.10479*.

Appendices

Within this supplementary material, we elaborate on the following aspects:

- Appendix A: Limitations and Future Work
- Appendix B: Safeguarding
- Appendix C: More Evaluation Details
- Appendix D: Prompt Details
- Appendix E: Causal Analysis in Details
- Appendix F: Dataset Details and License

A Limitations and Future Work

A.1 Limitations

Ethical and Applicability Considerations. U2-BENCH is designed as a research-oriented benchmark and is not intended for clinical deployment or diagnostic decision-making. Any real-world application of models evaluated on this benchmark would require separate validation and regulatory approval. Although all data sources are licensed or publicly available and de-identified where applicable, we acknowledge that not all ethical and demographic dimensions of fairness can be fully accounted for at this stage.

Evaluation Scope. The benchmark focuses on key task categories relevant to ultrasound interpretation—such as anatomical recognition, diagnostic classification, and structured report generation. While these tasks are representative and grounded in clinical utility, they do not exhaust the full landscape of sonographic applications. The evaluation metrics used (e.g., accuracy, BLEU) may not capture the full subtlety of expert clinical judgment, especially in edge cases.

Ultrasound-Specific Challenges. Ultrasound imaging is highly operator-dependent and subject to artifacts such as shadowing, speckle, and angle variation. Variability in scanning protocols and lack of standardized definitions (e.g., for "standard planes") can complicate model training and evaluation. These modality-specific challenges are inherent to ultrasound and reflect real-world complexities rather than flaws in the benchmark design.

A.2 Future Work

Extending Dataset Diversity and Robustness. While U2-BENCH aggregates data from a broad range of sources, further expansion to include more institutions, device types, and global populations would improve its representativeness. Future iterations of the benchmark will explore domain adaptation, adversarial robustness, and performance under distribution shifts to better simulate deployment conditions in varied clinical environments.

Model Generalization and Multimodal Reasoning. Current LVLMs still struggle with fine-grained spatial tasks, consistency across subgroups, and robust generation of clinically meaningful language. In future work, we aim to incorporate richer contextual information (e.g., patient history, multi-view inputs) to better assess models' multimodal integration capabilities and real-world reasoning performance.

Video-Based and Real-Time Evaluation. U2-BENCH currently operates on frame-based inputs to ensure comparability across models. However, clinical ultrasound interpretation often involves dynamic, probe-controlled acquisition. Extending the benchmark to include video sequences, real-time tasks, and longitudinal case studies will be a major step toward closing the simulation-to-clinic gap.

Theoretical Foundations and Causality. Our current benchmark is designed for practical performance evaluation. Future work will incorporate diagnostic reasoning audits, causal probing methods, and uncertainty quantification frameworks to deepen our understanding of LVLM behavior in high-stakes medical applications.

Standardization in Ultrasound AI. There is a growing need for community consensus on annotation standards, task definitions, and evaluation protocols in ultrasound AI. We hope U2-BENCH can serve as a starting point for these conversations and will actively evolve in response to feedback from both clinical and technical communities.

B Safeguarding

This study involves secondary use of de-identified, publicly available or licensed ultrasound datasets for the purpose of benchmarking machine learning models. All data used in **U2-BENCH** are either publicly released with appropriate usage permissions or obtained through official licensing agreements. No personally identifiable information is used, and all experiments are conducted in accordance with relevant data protection and ethical guidelines. Human annotators involved in quality assurance were trained to follow data confidentiality protocols, and no clinical decision-making was involved at any stage of this work.

C More Evaluation Details

C.1 Justification of U2-Score Weighting

The U2-Score summarizes model performance across the eight benchmark tasks in **U2-BENCH** through a weighted aggregation:

$$\text{U2-Score} := \sum_{t=1}^N w_t \cdot d_t, \quad \text{where} \quad w_t = \frac{n_t}{\sum_j n_j}, \quad d_t \in [0, 1] \quad (3)$$

Each task t is associated with a weight w_t proportional to its number of annotated examples n_t , and a normalized evaluation score d_t representing performance on that task. This formulation ensures that the final score reflects both task competence and dataset composition.

The weighting design of U2-Score is rooted in data-driven representation of ultrasound practice. All benchmark tasks are constructed from licensed and publicly available datasets sourced through a comprehensive and systematic search. As a result, the sample distribution across tasks captures the structure of real-world ultrasound data availability and usage. Tasks with greater sample counts typically correspond to more standardized, clinically widespread, and institutionally supported workflows—such as disease diagnosis and view recognition. Conversely, tasks with fewer examples often reflect more specialized applications or less frequently annotated modalities.

By aligning task weights with dataset size, the U2-Score encodes the practical significance and maturity of each task in the public ultrasound domain. This approach also promotes reliable aggregate evaluation, as scores are more robustly informed by tasks with greater data coverage. The resulting metric balances comprehensiveness with grounded applicability, supporting fair and reproducible model comparison in the current landscape of ultrasound AI research.

C.2 Error Analysis

Diagnosis Task 40: case001273

Accuracy: 0.05 (1/20 models correct)

Reference Answer: 4A

Model Responses:

- **gemini-2.5-pro-preview-03-25:** 5
- **Do-3B:** 2
- **gpt-4o-2024-08-06:** '3'
- **doubao-1.5-vision-pro-32k-250115:** Insufficient information provided to determine BI - RADS category. Please describe the sonographic characteristics of the lesion in the image.
- **MedDr:** 5
- **qwen-max-2025-01-25:** 3
- **local_model:** 2
- **claude-3-7-sonnet-20250219:** 4C
- **gemini-2.0-pro-exp-02-05:** 4C
- **deepseek-vl2:** 4B
- **qwen2.5-vl-3b-instruct:** 2
- **Qwen2.5-VL-7B-Instruct-Pro:** 4A
- **gemini-1.5-pro-latest:** 2
- **Qwen2.5-VL-32B-Instruct:** 4B
- **InternVL3-9B-Instruct:** 4B
- **Phi-4-multimodal-instruct:** I cannot synthesize the sonographic characteristics of the lesions or make a BI-RADS assessment without the actual ultrasound images or a detailed description of the findings. If you provide the relevant details or images, I would be able to assist you with the assessment.
- **LLaVA-1.5-13B-HF:** '2'
- **Mistral-Small-3.1-24B-Instruct-2503:** 3
- **Do-7B:** The anatomical component most vividly portrayed is the breast.
- **gpt-4o-mini:** 2

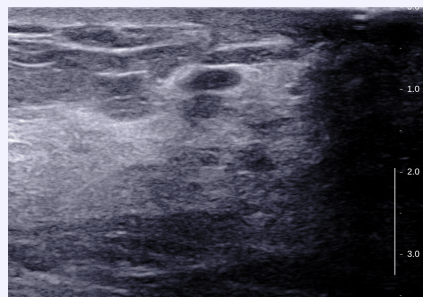


Figure 5: Ultrasound image for Diagnosis Task 40: case001273

Prompt: You are a radiologist analyzing a breast ultrasound image. Your task is to synthesize the sonographic characteristics of any identified lesions (or lack thereof) into a final ACR BI-RADS (Breast Imaging Reporting and Data System) assessment category.

BI-RADS Ultrasound Assessment Category Definitions

- '2' (Benign): Findings are definitively benign (e.g., simple cysts, intramammary lymph nodes, stable surgical implants/changes). 0% likelihood of malignancy. Requires routine screening follow-up.
 - '3' (Probably Benign): Findings have characteristic benign features but are not definitively benign (e.g., presumed fibroadenoma, complicated cyst). Very low likelihood of malignancy (<2%). Short-interval (e.g., 6-month) follow-up is typically recommended.
 - '4A' (Low Suspicion for Malignancy): Findings warrant biopsy but have a low probability of malignancy (>2% to ≤10%).
 - '4B' (Moderate Suspicion for Malignancy): Findings warrant biopsy with an intermediate probability of malignancy (>10% to ≤50%).
 - '4C' (High Suspicion for Malignancy): Findings warrant biopsy with a high probability of malignancy (>50% to <95%), without the classic features of Category 5.
 - '5' (Highly Suggestive of Malignancy): Findings have classic malignant features (e.g., irregular spiculated mass). Very high probability of malignancy (≥95%). Biopsy is required, and definitive action should be taken regardless of pathology results if discordant.
- Choose the single most appropriate BI-RADS assessment category from the options below.
options: ['2', '3', '4A', '4B', '4C', '5']
- Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

View Recognition and Assessment Tasks (VRA)



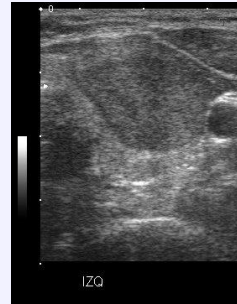
Reference Answer: hdvb

Model Responses:

- **Random Guessing:** huvf
- **MiniGPT-Med:** hdvf
- **MedDr:** hdvb
- **Qwen-2.5-VL-3B-Instruct:** hdvf
- **Qwen-2.5-VL-7B-Instruct:** hdvb
- **Qwen-2.5-VL-32B-Instruct:** hdvb
- **Qwen-2.5-VL-72B-Instruct:** hdvb
- **DeepSeek-VL2:** hdvf
- **InternVL3-9B-Instruct:** hdvf
- **LLaVA-1.5-13B:** huvb
- **Phi-4-Multimodal-Instruct:** hdvf
- **Mistral-Small-3.1-24B-Inst.:** hdvb
- **Doubao-1.5-Vision-Pro-32k:** hdvb
- **GPT-4o-Mini:** hdvf
- **GPT-4o:** hdvb
- **Gemini-1.5-Pro:** hdvf
- **Gemini-2.0-Pro-Exp:** hdvb
- **Gemini-2.5-Pro-Preview:** hdvf
- **Claude-3.7-Sonnet:** huvb
- **Qwen-Max:** hdvb
- **Dolphin-V1:** hdvb

Prompt: You are a radiologist analyzing a fetal ultrasound image. Your task is to determine the fetal presentation and orientation based on the provided ultrasound image. Specifically, identify whether the fetal head is down(hd) or up(hu). Additionally, determine if the fetal back is primarily oriented towards the ultrasound probe (vb) or towards the ultrasound probe (vf). Choose the single best option from the options below that accurately combines these findings. options: 'hdvb', 'hdvf', 'huvb', 'huvf' Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Lesion Localization Tasks (LL)



Reference Answer: upper left

Model Responses:

- **Random Guessing:** lower right
- **MiniGPT-Med:** upper center
- **MedDr:** upper left
- **Qwen-2.5-VL-3B-Instruct:** middle left
- **Qwen-2.5-VL-7B-Instruct:** upper left
- **Qwen-2.5-VL-32B-Instruct:** upper left
- **Qwen-2.5-VL-72B-Instruct:** upper left
- **DeepSeek-VL2:** upper right
- **InternVL3-9B-Instruct:** upper center
- **LLaVA-1.5-13B:** middle left
- **Phi-4-Multimodal-Instruct:** upper center
- **Mistral-Small-3.1-24B-Inst.:** upper left
- **Doubao-1.5-Vision-Pro-32k:** upper left
- **GPT-4o-Mini:** upper right
- **GPT-4o:** upper left
- **Gemini-1.5-Pro:** upper right
- **Gemini-2.0-Pro-Exp:** upper left
- **Gemini-2.5-Pro-Preview:** upper right
- **Claude-3.7-Sonnet:** upper center
- **Qwen-Max:** middle left
- **Dolphin-V1:** upper left

Prompt: You are a radiologist analyzing an ultrasound image of thyroid. Your task is to identify the primary location of any visible lesion(s) relative to the boundaries of the displayed image. Consider the lesion's center location or most prominent area when deciding. Choose the single option from the list below that best describes this location, even if the fit is approximate. options: upper left, upper center, upper right, middle left, center, middle right, lower left, lower center, lower right, not visible Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Organ Detection Tasks (OD)



Reference Answer: center

Model Responses:

- **Random Guessing:** lower left
- **MiniGPT-Med:** middle right
- **MedDr:** center
- **Qwen-2.5-VL-3B-Instruct:** middle right
- **Qwen-2.5-VL-7B-Instruct:** center
- **Qwen-2.5-VL-32B-Instruct:** center
- **Qwen-2.5-VL-72B-Instruct:** center
- **DeepSeek-VL2:** middle left
- **InternVL3-9B-Instruct:** lower center
- **LLaVA-1.5-13B:** middle right
- **Phi-4-Multimodal-Instruct:** middle right
- **Mistral-Small-3.1-24B-Inst.:** center
- **Doubao-1.5-Vision-Pro-32k:** center
- **GPT-4o-Mini:** lower center
- **GPT-4o:** center
- **Gemini-1.5-Pro:** lower center
- **Gemini-2.0-Pro-Exp:** center
- **Gemini-2.5-Pro-Preview:** lower center
- **Claude-3.7-Sonnet:** center
- **Qwen-Max:** middle right
- **Dolphin-V1:** center

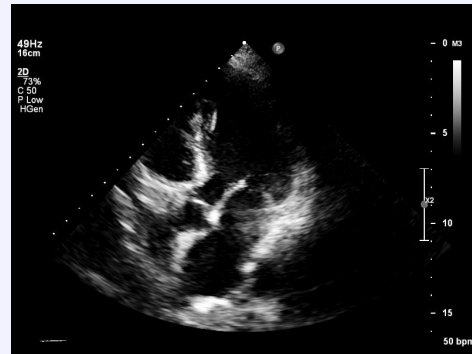
Prompt: You are a radiologist analyzing an ultrasound image of liver. Your task is to identify the primary location of the target organ relative to the boundaries of the displayed image. Consider the organ's center location or most prominent area when deciding. Choose the single option from the list below that best describes this location, even if the fit is approximate. options: upper left, upper center, upper right, middle left, center, middle right, lower left, lower center, lower right, not visible Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Keypoint Detection Tasks (KD)

Reference Answer: middle right

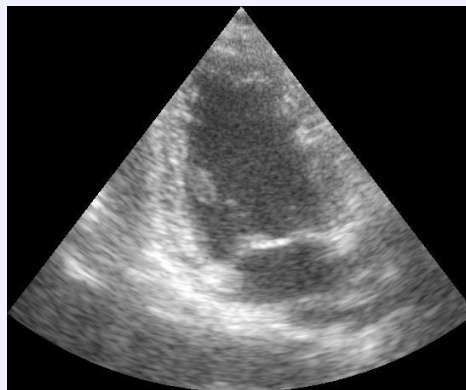
Model Responses:

- **Random Guessing:** upper center
- **MiniGPT-Med:** middle left
- **MedDr:** middle right
- **Qwen-2.5-VL-3B-Instruct:** center
- **Qwen-2.5-VL-7B-Instruct:** middle right
- **Qwen-2.5-VL-32B-Instruct:** middle right
- **Qwen-2.5-VL-72B-Instruct:** middle right
- **DeepSeek-VL2:** center
- **InternVL3-9B-Instruct:** middle left
- **LLaVA-1.5-13B:** center
- **Phi-4-Multimodal-Instruct:** lower right
- **Mistral-Small-3.1-24B-Inst.:** middle right
- **Doubao-1.5-Vision-Pro-32k:** middle right
- **GPT-4o-Mini:** center
- **GPT-4o:** middle right
- **Gemini-1.5-Pro:** center
- **Gemini-2.0-Pro-Exp:** middle right
- **Gemini-2.5-Pro-Preview:** center
- **Claude-3.7-Sonnet:** middle right
- **Qwen-Max:** center
- **Dolphin-V1:** middle right



Prompt: You are a radiologist analyzing an ultrasound image of heart. Your task is to identify the primary location of the key anatomical landmark point relative to the boundaries of the displayed image. Consider the landmark's precise position when deciding. Choose the single option from the list below that best describes this location, even if the fit is approximate. options: upper left, upper center, upper right, middle left, center, middle right, lower left, lower center, lower right, not visible Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Cardiac View Evaluation Tasks (CVE)



Reference Answer: 2CH

Model Responses:

- **Random Guessing:** 4CH
- **MiniGPT-Med:** 4CH
- **MedDr:** 2CH
- **Qwen-2.5-VL-3B-Instruct:** 4CH
- **Qwen-2.5-VL-7B-Instruct:** 4CH
- **Qwen-2.5-VL-32B-Instruct:** 4CH
- **Qwen-2.5-VL-72B-Instruct:** 4CH
- **DeepSeek-VL2:** 4CH
- **InternVL3-9B-Instruct:** 4CH
- **LLaVA-1.5-13B:** 4CH
- **Phi-4-Multimodal-Instruct:** 4CH
- **Mistral-Small-3.1-24B-Inst.:** 4CH
- **Doubao-1.5-Vision-Pro-32k:** 2CH
- **GPT-4o-Mini:** 4CH
- **GPT-4o:** 4CH
- **Gemini-1.5-Pro:** 4CH
- **Gemini-2.0-Pro-Exp:** 4CH
- **Gemini-2.5-Pro-Preview:** 4CH
- **Claude-3.7-Sonnet:** 2CH
- **Qwen-Max:** 4CH
- **Dolphin-V1:** 2CH

Prompt: You are a radiologist or cardiologist specializing in echocardiography, analyzing an apical view ultrasound image of the human heart.

Your task is to accurately identify the specific apical view presented in the provided echocardiogram image. Carefully examine the cardiac structures visible. Determine if the image displays primarily the left ventricle and left atrium only (indicative of a 2-Chamber view, 2CH), or if it clearly shows all four chambers: the left ventricle, right ventricle, left atrium, and right atrium (indicative of a 4-Chamber view, 4CH). Choose the single best option from the list below that correctly identifies the view.

options: 2CH, 4CH

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.



Reference Answer: Moderate OA

Model Responses:

- **Random Guessing:**
- **MiniGPT-Med:** Questionable OA
- **MedDr:** Moderate OA
- **Qwen-2.5-VL-3B-Instruct:** No OA
- **Qwen-2.5-VL-7B-Instruct:** Questionable OA
- **Qwen-2.5-VL-32B-Instruct:** Mild OA
- **Qwen-2.5-VL-72B-Instruct:** Mild OA
- **DeepSeek-VL2:** Questionable OA
- **InternVL3-9B-Instruct:** No OA
- **LLaVA-1.5-13B:** No OA
- **Phi-4-Multimodal-Instruct:** Questionable OA
- **Mistral-Small-3.1-24B-Inst.:** Mild OA
- **Doubao-1.5-Vision-Pro-32k:** Moderate OA
- **GPT-4o-Mini:** No OA
- **GPT-4o:** No OA
- **Gemini-1.5-Pro:** Questionable OA
- **Gemini-2.0-Pro-Exp:** Mild OA
- **Gemini-2.5-Pro-Preview:** Questionable OA
- **Claude-3.7-Sonnet:** Mild OA
- **Qwen-Max:** Mild OA
- **Dolphin-V1:** Moderate OA

Prompt: You are a radiologist analyzing an ultrasound image of left/right knee.

Your task is to assess the severity of osteoarthritis (OA) using the established Kellgren-Lawrence (KL) grading system. Kellgren-Lawrence (KL) Grade Mapping to Options:

- 'No OA': Corresponds to KL Grade 0 (No radiographic features of OA).
 - 'Questionable OA': Corresponds to KL Grade 1 (Doubtful JSN and possible minute osteophytes).
 - 'Mild OA': Corresponds to KL Grade 2 (Definite osteophytes and possible JSN).
 - 'Moderate OA': Corresponds to KL Grade 3 (Moderate multiple osteophytes, definite JSN, some sclerosis, possible deformity).
 - 'Severe OA': Corresponds to KL Grade 4 (Large osteophytes, marked JSN, severe sclerosis, definite deformity).
 - 'Total joint replacement': Indicates the presence of knee arthroplasty components (prosthesis), which replaces the native joint structures evaluated by the KL scale.
- options: 'Mild OA', 'Moderate OA', 'No OA', 'Questionable OA', 'Severe OA', 'Total joint replacement'

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Report Generation Tasks (RG) Input



Prompt: You are a radiologist analyzing an ultrasound image focused on the Liver. Your task is generate a concise and informative radiological report based strictly on the visual findings within the provided image. Your report should describe the primary organ's appearance (size, shape, borders/capsule), its parenchymal echotexture (e.g., homogeneous, heterogeneous, echogenicity relative to reference structures), and identify any visible abnormalities (e.g., masses, cysts, fluid collections, calcifications, ductal dilation). Comment on relevant adjacent structures if visualized. Use standard radiological terminology. Output format: Strings, that is your report.

Example: The liver morphology is full with a smooth capsule. The parenchymal echotexture is fine and diffusely increased. Visualization of the portal venous system is suboptimal. Intrahepatic and extrahepatic bile ducts are not dilated. The main portal vein diameter is within normal limits. The gallbladder is normal in size and shape. The wall is smooth and not thickened. No obvious abnormal echoes are seen within the lumen. The pancreas is normal in size and shape with homogeneous parenchymal echotexture. The pancreatic duct is not dilated. No definite space-occupying lesion is seen within the pancreas. The spleen is normal in size and shape with homogeneous parenchymal echotexture. No obvious space-occupying lesion is seen within the spleen.

Report Generation Tasks (RG) Output

Reference Answer: The liver morphology is full with a smooth capsule. The parenchymal echotexture is fine and diffusely increased. Visualization of the portal venous system is suboptimal. Intrahepatic and extrahepatic bile ducts are not dilated. The main portal vein diameter is within normal limits. No focal lesions are identified.

Model Responses:

- **Random Guessing:**
- **MiniGPT-Med:** The liver shows normal size and contour with homogeneous echotexture. No focal lesions or biliary dilatation noted. Portal structures appear normal.
- **MedDr:** The liver demonstrates normal size and contour with a smooth capsule. The parenchymal echotexture is homogeneous with normal echogenicity. No focal lesions identified. Intrahepatic biliary ducts are not dilated. Portal venous structures appear normal.
- **Qwen-2.5-VL-3B-Instruct:** The liver appears normal in size with a smooth contour. The parenchymal echotexture is homogeneous without focal lesions. No intrahepatic biliary dilatation is seen.
- **Qwen-2.5-VL-7B-Instruct:** The liver appears normal in size with a smooth capsule. The parenchymal echotexture is homogeneous with no focal lesions. No intrahepatic or extrahepatic biliary dilatation is seen. The portal venous structures appear normal.
- **Qwen-2.5-VL-32B-Instruct:** The liver appears normal in size with a smooth contour. The parenchymal echotexture is homogeneous with normal echogenicity. No focal lesions such as masses, cysts, or calcifications are identified. The intrahepatic bile ducts are not dilated. The portal venous structures appear normal. No perihepatic fluid collection is observed.
- **Qwen-2.5-VL-72B-Instruct:** The liver demonstrates normal size and contour with a smooth capsule. The parenchymal echotexture is homogeneous with normal echogenicity. No focal lesions or masses are identified. The intrahepatic biliary tree is not dilated. Portal venous structures appear normal. No perihepatic fluid collection is seen.
- **DeepSeek-VL2:** The liver appears normal in size with a smooth contour. The parenchymal echotexture is homogeneous with normal echogenicity. No focal lesions are identified. The intrahepatic bile ducts are not dilated. The portal venous structures appear normal. No perihepatic fluid collection is observed.
- **InternVL3-9B-Instruct:** The liver appears normal in size with a smooth contour. The parenchymal echotexture is homogeneous without focal lesions. No intrahepatic or extrahepatic biliary dilatation is seen. The portal venous structures appear normal.
- **LLaVA-1.5-13B:** The liver appears normal in size with a smooth contour. The parenchymal echotexture is homogeneous without focal lesions. No intrahepatic biliary dilatation is seen. The portal structures appear normal.
- **Phi-4-Multimodal-Instruct:** The liver appears normal in size with a smooth contour. The parenchymal echotexture is homogeneous with no focal lesions. No intrahepatic or extrahepatic biliary dilatation is seen. The hepatic vessels appear normal.
- **Mistral-Small-3.1-24B-Inst.:** The liver demonstrates normal size and contour with a smooth capsule. The parenchymal echotexture is homogeneous with normal echogenicity. No focal lesions or masses are identified. The intrahepatic biliary tree is not dilated. Portal venous structures appear normal.
- **Doubao-1.5-Vision-Pro-32k:** The liver appears normal in size with a smooth contour. The parenchymal echotexture is homogeneous with normal echogenicity. No focal lesions are identified. The intrahepatic bile ducts are not dilated. The portal venous structures appear normal.
- **GPT-4o-Mini:** The liver appears normal in size with a smooth contour. The parenchymal echotexture is homogeneous without focal lesions. No intrahepatic or extrahepatic biliary dilatation is seen. The hepatic vessels appear normal. No ascites or other abnormalities are identified in the visualized portion of the liver.

- **GPT-4o:** The liver appears normal in size and shape with a smooth capsule. The parenchymal echotexture is homogeneous with no obvious focal lesions. No intrahepatic or extrahepatic bile duct dilation is noted. The gallbladder is not clearly visualized in this image. No ascites or other fluid collections are present. Adjacent structures are not well visualized for further evaluation.
- **Gemini-1.5-Pro:** The liver appears normal in size with a smooth contour. The parenchymal echotexture is homogeneous with normal echogenicity. No focal lesions such as masses, cysts, or calcifications are identified. The intrahepatic bile ducts are not dilated. The portal venous structures appear normal.
- **Gemini-2.0-Pro-Exp:** The liver appears normal in size with a smooth capsule. The parenchymal echotexture is homogeneous with normal echogenicity. No focal lesions are identified. The intrahepatic bile ducts are not dilated. The portal venous structures appear normal where visualized.
- **Gemini-2.5-Pro-Preview:** The liver appears normal in size with a smooth capsule. The parenchymal echotexture is homogeneous with normal echogenicity. No focal lesions such as masses, cysts, or calcifications are identified. The intrahepatic bile ducts are not dilated. The portal venous structures appear normal where visualized. No perihepatic fluid collection is seen. The gallbladder and other adjacent structures are not well visualized in this limited image. Overall impression: Normal liver parenchyma without evidence of focal lesions or biliary dilation.
- **Claude-3.7-Sonnet:** The liver demonstrates normal size and contour with a smooth capsule. The parenchymal echotexture appears homogeneous with mildly increased echogenicity, suggesting possible early fatty infiltration. No focal lesions, masses, or cysts are identified. The intrahepatic biliary tree is not dilated. Portal venous structures appear normal where visualized. No perihepatic fluid collection is seen. Adjacent structures are not well visualized in this limited view.
- **Qwen-Max:** The liver appears normal in size with a smooth contour. The parenchymal echotexture is homogeneous with normal echogenicity. No focal lesions such as masses, cysts, or calcifications are identified. The intrahepatic bile ducts are not dilated. The portal venous structures appear normal. No perihepatic fluid collection is observed. The visualized portions of adjacent structures appear unremarkable. Overall, the liver appears sonographically normal based on this limited view.
- **Dolphin-V1:** The liver morphology is full with a smooth capsule. The parenchymal echotexture is fine and diffusely increased. Visualization of the portal venous system is suboptimal. Intrahepatic and extrahepatic bile ducts are not dilated. The main portal vein diameter is within normal limits. No focal lesions are identified.

To provide deeper insights into model performance on medical imaging tasks, we conduct a comprehensive error analysis of models across four critical ultrasound image analysis tasks: measurement, classification, segmentation, and report generation. This analysis reveals distinct error patterns and task-specific challenges that inform future model improvements.

Numerical Regression Task Analysis Among 101 total responses, the most significant challenge is the prevalence of **descriptive responses instead of numerical values** (53.47%). The model frequently generates interpretative text such as “The principal anatomical element visualized here is unequivocally the fetus head” rather than the expected numerical measurement (e.g., 291.4mm). This pattern suggests fundamental misunderstanding of task requirements, where the model interprets the task as image content identification rather than quantitative measurement.

Format violations constitute 1.98% of responses, where models provide numerical values with units (e.g., “113.6 mm”) despite explicit formatting constraints. Notably, 43.56% of responses follow the correct numerical format, though accuracy assessment requires comparison with ground truth values. The high rate of descriptive responses indicates that current vision-language models struggle with the transition from visual analysis to precise quantitative output.

Classification Task Performance Classification tasks demonstrate superior format compliance compared to measurement tasks, with 75.66% of responses providing valid option selections from 152 total responses. However, two distinct error patterns emerge: **explanatory responses** (5.92%) where models provide justifications rather than selections (e.g., “There is no definitive view of the fetal abdomen or pelvis to determine fetal position”), and **format violations** (18.42%) containing additional descriptive content alongside valid options.

The tendency toward explanatory responses reveals an interesting model behavior where excessive caution leads to task avoidance rather than best-effort selection from available options. This suggests that models may benefit from more explicit instructions emphasizing the requirement for definitive option selection even under uncertainty.

Segmentation and Localization Analysis Segmentation tasks, requiring spatial reasoning for anatomical structure localization, show moderate success with 66% valid position responses from 500 total responses. The primary error categories include **invalid position terminology** (27.80%) with responses like “Not visible.” or “Upper right.” that contain punctuation or non-standard terms, and **complete task deviation** (6.20%) where models provide structural descriptions instead of positional information.

Case Study Examples: Analysis of specific segmentation cases reveals distinct model behaviors. In thyroid lesion localization tasks, while Gemini-2.5-Pro and GPT-4o consistently provide concise responses (“center”), Claude-3.7-Sonnet exhibits significant format violations. For instance, when tasked with identifying tumor location in breast ultrasound images, Claude generated extensive explanatory text:

“This image appears to be an ultrasound showing tissue layers with varying echogenicity... I cannot identify a clear, definitive lesion... For proper medical diagnosis, this ultrasound should be evaluated by a qualified radiologist...”

Such responses, while demonstrating medical awareness, completely violate the specified output format requiring only location terms. This pattern suggests that Claude prioritizes safety disclaimers over task compliance in medical contexts.

Additionally, a concerning pattern emerges where multiple models consistently respond “center” regardless of actual lesion position, as evidenced by reference bounding boxes indicating lesions at coordinates [0.6, 0.247] and [0.595, 0.308]. This suggests potential spatial reasoning limitations or default response bias that could compromise clinical utility.

The relatively high success rate in spatial localization compared to numerical measurement suggests that discrete spatial reasoning may be more accessible to current vision-language architectures than continuous numerical estimation.

Report Generation Excellence Report generation tasks achieve the highest success rate (98%) among all evaluated tasks, with only 2% exhibiting structural misidentification and 1% showing false

findings. The rare but critical errors include anatomical misidentification (“Top view of fetus head and thorax” for fetal head ultrasound) and false pathological findings (“Aneuploid fetus with abnormal facial features”). While infrequent, such errors carry significant clinical implications, potentially leading to unnecessary medical interventions or patient anxiety.

Cross-Task Error Pattern Analysis Task difficulty ranking from most to least challenging reveals: measurement (43.56% success) > segmentation (66% success) > classification (75.66% success) > report generation (98% success). This hierarchy reflects the increasing complexity of transitioning from free-form text generation to structured, constrained outputs requiring precise adherence to format specifications.

Common error patterns across tasks include: (1) **descriptive language substitution**, most prominent in measurement tasks where models default to interpretative text rather than required numerical values; (2) **format non-compliance**, prevalent across classification and segmentation tasks despite clear formatting instructions; and (3) **task misunderstanding**, where models completely misinterpret task objectives, such as treating localization as structure identification.

Implications for Medical AI Development These findings highlight critical considerations for deploying vision-language models in medical imaging applications. The inverse relationship between task constraint and model performance suggests that current architectures excel at unconstrained text generation but struggle with precise, structured outputs essential for clinical decision-making. Future developments should prioritize: (1) enhanced instruction following capabilities for constrained output generation, (2) domain-specific fine-tuning on medical imaging tasks emphasizing numerical precision, and (3) robust validation mechanisms to detect and prevent false findings in clinical applications.

The analysis underscores that while large vision-language models show promise for medical imaging applications, careful task-specific optimization and human oversight remain essential, particularly for quantitative measurements and diagnostic assessments where precision directly impacts patient care.

D Prompt for Tasks

Prompt Template used for fetal view classification (dataset 10)

You are a radiologist analyzing a fetal ultrasound image.

Your task is to determine the fetal presentation and orientation based on the provided ultrasound image. Specifically, identify whether the fetal head is down(hd) or up(hu). Additionally, determine if the fetal back is primarily oriented towards the ultrasound probe (vb) or towards the ultrasound probe (vf). Choose the single best option from the options below that accurately combines these findings.

options: 'hdvb', 'hdvf', 'huvb', 'huvf'

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for heart view classification (dataset 18)

You are a radiologist or cardiologist specializing in echocardiography, analyzing an apical view ultrasound image of the human heart.

Your task is to accurately identify the specific apical view presented in the provided echocardiogram image. Carefully examine the cardiac structures visible. Determine if the image displays primarily the left ventricle and left atrium only (indicative of a 2-Chamber view, 2CH), or if it clearly shows all four chambers: the left ventricle, right ventricle, left atrium, and right atrium (indicative of a 4-Chamber view, 4CH). Choose the single best option from the list below that correctly identifies the view.

options: 2CH, 4CH

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for (KL) grading (dataset 28)

You are a radiologist analyzing an ultrasound image of left/right knee.

Your task is to assess the severity of osteoarthritis (OA) using the established Kellgren-Lawrence (KL) grading system. Kellgren-Lawrence (KL) Grade Mapping to Options:

- 'No OA': Corresponds to KL Grade 0 (No radiographic features of OA).
- 'Questionable OA': Corresponds to KL Grade 1 (Doubtful JSN and possible minute osteophytes).
- 'Mild OA': Corresponds to KL Grade 2 (Definite osteophytes and possible JSN).
- 'Moderate OA': Corresponds to KL Grade 3 (Moderate multiple osteophytes, definite JSN, some sclerosis, possible deformity).
- 'Severe OA': Corresponds to KL Grade 4 (Large osteophytes, marked JSN, severe sclerosis, definite deformity).
- 'Total joint replacement': Indicates the presence of knee arthroplasty components (prosthesis), which replaces the native joint structures evaluated by the KL scale.

Choose the single best option from the following list that accurately describes the image.

options: 'Mild OA', 'Moderate OA', 'No OA', 'Questionable OA', 'Severe OA', 'Total joint replacement'

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for BI-RADS classification (dataset 40)

You are a radiologist analyzing a breast ultrasound image.

Your task is to synthesize the sonographic characteristics of any identified lesions (or lack thereof) into a final ACR BI-RADS (Breast Imaging Reporting and Data System) assessment category.

BI-RADS Ultrasound Assessment Category Definitions:

- '2' (Benign): Findings are definitively benign (e.g., simple cysts, intramammary lymph nodes, stable surgical implants/changes). 0% likelihood of malignancy. Requires routine screening follow-up.
- '3' (Probably Benign): Findings have characteristic benign features but are not definitively benign (e.g., presumed fibroadenoma, complicated cyst). Very low likelihood of malignancy (<2%). Short-interval (e.g., 6-month) follow-up is typically recommended.
- '4A' (Low Suspicion for Malignancy): Findings warrant biopsy but have a low probability of malignancy (>2% to ≤10%).
- '4B' (Moderate Suspicion for Malignancy): Findings warrant biopsy with an intermediate probability of malignancy (>10% to ≥50%).
- '4C' (High Suspicion for Malignancy): Findings warrant biopsy with a high probability of malignancy (>50% to <95%), without the classic features of Category 5.
- '5' (Highly Suggestive of Malignancy): Findings have classic malignant features (e.g., irregular spiculated mass). Very high probability of malignancy (≥95%). Biopsy is required, and definitive action should be taken regardless of pathology results if discordant.

Choose the single most appropriate BI-RADS assessment category from the options below.

options: ['2', '3', '4A', '4B', '4C', '5']

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for fetal abdomen (dataset 50)

You are a radiologist analyzing an ultrasound image of fetal abdomen.

Your task is to determine if the presented cross-sectional view of the fetal abdomen is technically adequate for performing an accurate Abdominal Circumference (AC) measurement according to standard obstetric guidelines. Identify the specific anatomical plane shown for the fetal abdomen. Determine if this plane meets the criteria for an optimal AC measurement (correct landmarks visible, proper transverse orientation) or if it is suboptimal (incorrect plane, missing landmarks, oblique/foreshortened view, presence of interfering structures like kidneys). Choose the single best option describing the plane's suitability for AC measurement.

options: 'none', 'optimal', 'suboptimal'

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for breast classification

You are a radiologist analyzing a breast ultrasound image.

Your task is carefully examine the provided breast ultrasound image, evaluate any identified lesions or abnormalities based on key sonographic characteristics (including shape, orientation, margin, echo pattern, posterior acoustic features, and associated features), synthesize these features to form an overall impression about the likelihood of malignancy, and then choose the single best option from the following list that accurately summarizes this assessment.

options: (normal), benign, malignant

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for thyroid classification

You are a radiologist specializing in head and neck or endocrine imaging, analyzing an ultrasound image of the thyroid gland.

Your task is to carefully examine the provided thyroid ultrasound image, evaluate the overall thyroid gland parenchyma (echogenicity, texture, vascularity), identify any focal nodules, assess the specific sonographic features of any nodules found (including composition, echogenicity, shape, margin, and echogenic foci), synthesize these findings to determine if the gland appears normal, contains benign-appearing findings, or contains findings suspicious for malignancy, and then choose the single best option from the following list that accurately summarizes this assessment.

options: (normal thyroid), benign, malignant

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for skin cancer classification (dataset 25)

You are a radiologist analyzing an ultrasound image of skin.

Your task is to carefully examine the provided skin ultrasound image, evaluate the identified lesion or abnormality based on key sonographic characteristics (including its location within skin layers, echogenicity, internal echo texture, shape, margins, size/depth, posterior acoustic phenomena, and vascularity assessed with Doppler), synthesize these features to form an overall impression regarding the likelihood of malignancy, and then choose the single best option from the following list that summarizes this assessment.

options: benign, malignant

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for pancreas cancer classification (dataset 42)

You are a radiologist analyzing an ultrasound image of the pancreas.

Your task is to carefully examine the provided ultrasound image of the pancreas, evaluate the gland's echotexture, size, margins, and the pancreatic duct diameter, identify any focal lesions or masses (noting their echogenicity, margins, size, and vascularity if Doppler is available), assess for associated findings such as ductal dilation (including potential "double duct" sign), vascular involvement (encasement/thrombosis), regional lymphadenopathy, or fluid collections, synthesize these findings to determine if there is evidence suspicious for primary pancreatic cancer versus other findings, and then choose the single best option from the following list that summarizes this assessment.

options: non-pancreas cancer, pancreas cancer

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for PCOS classification (dataset 74)

You are a radiologist analyzing an ultrasound image obtained during a pelvic examination, potentially as part of an evaluation for Polycystic Ovary Syndrome (PCOS).

Your task is to evaluate the overall appearance of the anatomical structures presented in the ultrasound image (primarily focusing on the ovaries and potentially the uterus). Consider sonographic features such as ovarian size, morphology, follicle count and distribution, stromal echogenicity, as well as any other findings that might indicate pathology. Based on this assessment, determine if the image appears generally normal or if it displays features suggestive of an abnormality (which could include findings consistent with PCOS or other conditions). Choose the single best option from the following list that accurately describes this overall impression.

options: 'Appears abnormal', 'Appears normal'

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for PCOS classification (dataset 74)

You are a radiologist analyzing an ultrasound image obtained during a pelvic examination. Crucially, assume this specific image has already been determined to show some form of abnormality. Your focus now is on the nature of that abnormality.

Your task is to specifically assess whether the abnormality present in this ultrasound image includes clear sonographic evidence consistent with a polycystic ovary. Evaluate the visualized ovarian structures, paying close attention to features commonly associated with PCOS, such as: increased number of follicles, peripheral distribution of follicles, increased ovarian volume, increased stromal echogenicity or volume. Based on whether these specific PCOS-related sonographic features are identifiable within the overall abnormal appearance, specifies whether the ultrasound image shows evidence/ visibility of a polycystic ovary or not. Choose the single best option from the following list.

options: 'Not-visible', 'Visible'

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for PCOS classification (dataset 75)

You are a radiologist analyzing an ultrasound image obtained during a pelvic examination, specifically being evaluated for features potentially related to Polycystic Ovary Syndrome (PCOS).

Your task is to carefully evaluate the provided ultrasound image for sonographic features consistent with Polycystic Ovarian Morphology (PCOM), which is the ultrasound component relevant to PCOS detection. Analyze the visualized ovary (or ovaries), considering criteria such as increased ovarian volume, increased antral follicle count (e.g., ≥ 20 per ovary), peripheral follicle distribution, and / or increased stromal echogenicity / volume. If sonographic features consistent with PCOM are present, select the label 'infected', otherwise 'noninfected'. Choose the single best option from the following list.

options: 'infected', 'noninfected'

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for lung parenchyma (dataset 44)

You are a radiologist or clinician skilled in performing and interpreting Lung Ultrasound (LUS), specifically analyzing an ultrasound image of the lung pleura and parenchyma.

Your task is to carefully examine the provided lung ultrasound image, focusing on the appearance of the pleural line and the underlying lung parenchyma, identify the presence and characteristics of A-lines, B-lines (number, coalescence), and any consolidations according to the defined severity scoring criteria below, and then choose the single best integer score (0, 1, 2, or 3) from the following list that accurately reflects the observed findings.

LUS Severity Score Criteria:

- 0: Normal lung pattern. Characterized by a continuous, regular, thin pleural line with horizontal reverberation artifacts (A-lines) below it. Sliding lung sign is typically present.
- 1: Mild interstitial syndrome. Characterized by an indented or slightly irregular pleural line. Scattered, well-defined vertical artifacts (B-lines) are visible (typically ≥ 3 B-lines per intercostal space but not coalescent).
- 2: Moderate interstitial syndrome or early consolidation. Characterized by a broken or significantly irregular pleural line. Multiple coalescent B-lines (small "white lung" areas) or small subpleural consolidations are present.
- 3: Severe interstitial syndrome or large consolidation. Characterized by dense and largely extended confluent B-lines ("white lung" appearance occupying most or all of the screen) with or without large consolidations.

Options: 0, 1, 2, 3

Output format: only the single chosen integer number from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for fatty liver classification (dataset 57)

You are a radiologist analyzing a static B-mode ultrasound image displaying the liver.

Your task is to evaluate the liver parenchyma in the provided image to determine the grade of hepatic steatosis. For this task, label 1 is assigned if the image displays features consistent with fatty liver (which often correlates histologically with >5% hepatocyte steatosis), while label 0 is assigned if such features are absent. Based on your comprehensive assessment of these sonographic features, determine whether the image displays sufficient evidence to be classified as showing fatty liver (Label 1) or not (Label 0). Choose the single best option from the following list that accurately reflects your classification.

options: 0, 1

Output format: only the single chosen integer number from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for fetal (dataset 03)

You are a radiologist analyzing a single ultrasound image acquired during a fetal examination.

Your task is to carefully examine the provided image, identify the primary anatomical structure or region being visualized, and determine the most appropriate description based on the standard imaging planes used in fetal ultrasound. Choose the single best option from the following list that accurately describes the main subject shown in the image.

options: 'fetal abdomen', 'fetal femur', 'fetal brain', 'fetal thorax', 'maternal cervix', 'other'

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for thyroid plane classification (dataset 37)

You are a radiologist with expertise in interpreting neck and thyroid ultrasound images. You are presented with a single B-mode ultrasound image focused on the thyroid gland and adjacent neck structures.

Your task is to identify the Cardinal Anatomical Plane depicted in the provided ultrasound image. Choose the single best option from the following list that accurately describes the image.

options: 'Axial/Transverse Plane', 'Coronal Plane', 'Sagittal Plane'

Output format: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or formatting.

Prompt Template used for fetal (dataset 53)

You are a radiologist analyzing a single B-mode ultrasound image obtained during a fetal assessment.

Your task is to carefully examine the provided ultrasound image frame to identify the presence or absence of two specific anatomical landmarks: the fetal head and the maternal symphysis pubis. Based on this identification, classify the frame's content by choosing the single best option from the following list that accurately describes which of these landmarks are visible. Choose the single best option from the following list that accurately describes the frame's content.

options: 'None', 'OnlyFetalHead', 'OnlySymphysisPubis', 'SymphysisPubis+FetalHead'

Output prompt: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or other formatting.

Prompt Template used for carotid classification (dataset 69)

You are a radiologist analyzing an ultrasound image depicting a portion of the carotid arterial system in the neck.

Your task is to carefully examine the provided ultrasound image, analyzing anatomical landmarks, vessel morphology, and its position relative to other neck structures, to identify the primary carotid artery segment shown. Choose the single best option from the following list that accurately describes the main vessel visualized in the frame's content. Assume 'left carotid' and 'right carotid' refer generally to the common or internal carotid artery on that respective side, while 'external carotid' refers specifically to the external carotid artery branch. Choose the single best option from the following list that accurately describes the image.

options: 'external carotid', 'left carotid', 'right carotid'

Output prompt: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or other formatting.

Prompt Template used for anatomy classification

You are an expert specialized in analyzing medical ultrasound images. You are provided with a single ultrasound image frame, which could depict various parts of the human body.

Your task is to analyze the provided ultrasound image and identify the primary anatomical region or organ system being visualized. Choose the single best option from the following list that most accurately represents this primary anatomical subject.

options: 'fetal', 'thyroid', 'heart', 'lung', 'liver', 'carotid', 'kidney', 'prostate', 'breast', 'other'

Output prompt: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or other formatting.

Prompt Template used for knee classification

You are a radiologist analyzing an ultrasound image of knee.

Your task is to classify the specific anatomical view, laterality (left/right), orientation, and any specific imaging technique or patient positioning shown in the image:

- 'left anterior suprapatellar longitudinal': Image of the left knee, taken from the front (anterior), just above the kneecap (suprapatellar), with the ultrasound probe oriented along the long axis of the thigh/patellar tendon. Standard B-mode imaging.
- 'left anterior suprapatellar longitudinal with power Doppler': Same view as above (left, anterior suprapatellar, longitudinal), but with Power Doppler mode activated, typically used to assess blood flow or inflammation.
- 'left anterior suprapatellar transverse in 30 degrees flexion': Image of the left knee, from the front (anterior), above the kneecap (suprapatellar), with the probe oriented across (transverse) the thigh, and the knee bent at approximately 30 degrees.
- 'left anterior suprapatellar transverse in maximal flexion': Same view as above (left, anterior suprapatellar, transverse), but with the knee bent as much as possible (maximal flexion).
- 'left lateral longitudinal': Image of the outer side (lateral) of the left knee, with the probe oriented along the long axis of the structures (e.g., LCL, IT band).
- 'left medial longitudinal': Image of the inner side (medial) of the left knee, with the probe oriented along the long axis of the structures (e.g., MCL, medial meniscus).
- 'left posterior medial transverse': Image of the back, inner corner (posterior medial) of the left knee, with the probe oriented across (transverse) the structures (often used for Baker's cysts).
- 'right anterior suprapatellar longitudinal': Image of the right knee, taken from the front (anterior), just above the kneecap (suprapatellar), with the ultrasound probe oriented along the long axis of the thigh/patellar tendon. Standard B-mode imaging.
- 'right anterior suprapatellar longitudinal with power Doppler': Same view as above (right, anterior suprapatellar, longitudinal), but with Power Doppler mode activated.
- 'right anterior suprapatellar transverse in 30 degrees flexion': Image of the right knee, from the front (anterior), above the kneecap (suprapatellar), with the probe oriented across (transverse) the thigh, and the knee bent at approximately 30 degrees.
- 'right anterior suprapatellar transverse in maximal flexion': Same view as above (right, anterior suprapatellar, transverse), but with the knee bent as much as possible (maximal flexion).
- 'right lateral longitudinal': Image of the outer side (lateral) of the right knee, with the probe oriented along the long axis of the structures.
- 'right medial longitudinal': Image of the inner side (medial) of the right knee, with the probe oriented along the long axis of the structures.
- 'right posterior medial transverse': Image of the back, inner corner (posterior medial) of the right knee, with the probe oriented across (transverse) the structures.

Choose the single best option from the following list that accurately describes the image.

Options: 'left anterior suprapatellar longitudinal', 'left anterior suprapatellar longitudinal with power Doppler', 'left anterior suprapatellar transverse in 30 degrees flexion', 'left anterior suprapatellar transverse in maximal flexion', 'left lateral longitudinal', 'left medial longitudinal', 'left posterior medial transverse', 'right anterior suprapatellar longitudinal', 'right anterior suprapatellar longitudinal with power Doppler', 'right anterior suprapatellar transverse in 30 degrees flexion', 'right anterior suprapatellar transverse in maximal flexion', 'right lateral longitudinal', 'right medial longitudinal', 'right posterior medial transverse'

Output prompt: only the exact text of the chosen option from the list above. Do not include any introductory phrases, explanations, numbering, or other formatting.

Prompt Template used for lesion detection

You are a radiologist analyzing an ultrasound image of thyroid.

Your task is to identify the primary location of any visible lesion(s) relative to the boundaries of the displayed image. Consider the lesion's center location or most prominent area when deciding. Choose the single option from the list below that best describes this location, even if the fit is approximate.

Choose the single most appropriate location from the following list:

- upper left
- upper center
- upper right
- middle left
- center
- middle right
- lower left
- lower center
- lower right
- not visible

Output format: only one or two word(s) representing the chosen location. No additional text or formatting is allowed.

Prompt Template used for organ detection

You are a radiologist analyzing an ultrasound image of abdominal.

Your task is to determine the primary location, relative to the image boundaries, for each visible structure listed in liver.

- Consider the structure's center or most prominent area when deciding its location.
- Choose the single option from the list below that best describes the location, even if the fit is approximate.

Location Options:

- upper left
- upper center
- upper right
- middle left
- center
- middle right
- lower left
- lower center
- lower right
- not visible

Output format: only one or two word(s) representing the chosen location. No additional text or formatting is allowed.

Prompt Template used for keypoint detection

You are a radiologist analyzing an ultrasound image of the heart.

Your task is to determine the top inner point of the aortic valve.

- Consider the structure's center or most prominent area when deciding its location.
- Choose the single option from the list below that best describes the location, even if the fit is approximate.

Location Options:

- upper left
- upper center
- upper right
- middle left
- center
- middle right
- lower left
- lower center
- lower right
- not visible

Output format: only one or two word(s) representing the chosen location. No additional text or formatting is allowed.

Prompt Template used for caption generation

You are a radiologist analyzing an ultrasound image focused on the {anatomy_location}.

Your task is to generate a concise and informative caption that accurately describes the key anatomical structures and any significant findings visible in the provided ultrasound image.

Output format: A single string constituting the image caption. Output only the generated caption text itself. Do not include any introductory phrases (like "Caption:"), labels, explanations, or additional formatting.

Examples:

Example1: Thyroid nodule in the right lobe. TI-RADS level 3, Benign.

Example2: Thyroid nodule in the left lobe. TI-RADS level 3, Benign.

Example3: Thyroid nodule in the right lobe. TI-RADS level 4, Benign.

Prompt Template used for report generation

You are a radiologist analyzing an ultrasound image focused on the {anatomy_location}.

Your task is generate a concise and informative radiological report based strictly on the visual findings within the provided image. Your report should describe the primary organ's appearance (size, shape, borders/capsule), its parenchymal echotexture (e.g., homogeneous, heterogeneous, echogenicity relative to reference structures), and identify any visible abnormalities (e.g., masses, cysts, fluid collections, calcifications, ductal dilation). Comment on relevant adjacent structures if visualized. Use standard radiological terminology.

Output format: Strings, that is your report.

Example: The liver morphology is full with a smooth capsule. The parenchymal echotexture is fine and diffusely increased. Visualization of the portal venous system is suboptimal. Intrahepatic and extrahepatic bile ducts are not dilated. The main portal vein diameter is within normal limits. The gallbladder is normal in size and shape. The wall is smooth and not thickened. No obvious abnormal echoes are seen within the lumen. The pancreas is normal in size and shape with homogeneous parenchymal echotexture. The pancreatic duct is not dilated. No definite space-occupying lesion is seen within the pancreas. The spleen is normal in size and shape with homogeneous parenchymal echotexture. No obvious space-occupying lesion is seen within the spleen.

E Causal Analysis in Details

In this appendix, we provide a detailed causal interpretation of our prompt ablation experiment, based on structural causal modeling and informed by recent advances in causal prompting methods [86] for large language models (LLMs).

Structural Causal Model. Let X denote the prompt formulation (with or without anatomy), A the model’s output (e.g., prediction correctness), and U an unobserved confounder encapsulating latent model biases or corpus priors. We assume a structural causal model (SCM) with the graph $X \leftarrow U \rightarrow A$ and $X \rightarrow A$, indicating that the observed association between X and A may be confounded by U .

Causal Estimation via Front-Door Adjustment. Since U is unobserved and cannot be directly conditioned on, traditional back-door adjustment is infeasible. However, the causal prompting framework suggests that under front-door conditions, we can still estimate the causal effect of X on A by conditioning on an observed mediator R that lies on the causal path $X \rightarrow R \rightarrow A$ [86].

In our case, we do not explicitly use a chain-of-thought as a mediator, but we achieve equivalent control via a *paired evaluation design*, where each input image x_i is processed under both prompt conditions ($P_2 = 1$) and ($P_2 = 0$). Since the image and model remain unchanged across conditions, this implicitly blocks the back-door path through U , allowing us to estimate the interventional effect $P(A \mid \text{do}(P_2))$ as:

$$\mathbb{E}[A_1 - A_0] \approx P(A \mid \text{do}(P_2 = 1)) - P(A \mid \text{do}(P_2 = 0)) \quad (4)$$

Justifying the Ignorability of the Confounder. Our analysis relies on the key assumption that the confounder U is shared across paired samples. That is, any latent bias in the model remains fixed for a given input x_i regardless of the prompt variant. This symmetry mirrors the assumptions made in front-door adjustment where the mediator R is used to block paths from U to A [86]. By controlling X while holding U constant through sample pairing, we achieve a quasi-interventional setting:

$$A_1 - A_0 = f(P_2 = 1, U) - f(P_2 = 0, U) \quad (5)$$

This allows the confounding effect of U to cancel out.

McNemar’s Test for Effect Significance. To test whether the effect of including anatomical tokens is statistically significant, we apply McNemar’s test to paired binary outcomes across 521 samples. This evaluates the asymmetry of correct predictions between the two prompt settings. A significant χ^2 statistic (16.04, $p < 10^{-4}$) confirms a causal relationship between X and A .

Conclusion. Inspired by the front-door adjustment framework, our paired evaluation design provides a valid estimation of the causal effect of prompt modification, without requiring access to the confounder U . This supports the claim that including anatomy information in the prompt has a positive causal effect on model accuracy in ultrasound understanding.

F Dataset Details and License

Table 6: Summary of Annotated Datasets Used in U2-BENCH

| Dataset | Anatomy | Clinical scenarios | Task | Case | License |
|--|---|--|-------------------|------|-----------------|
| FETAL PLANES DB [12] | Fetal abdomen Fetal brain Fetal femur Fetal thorax Maternal cervix other | Fetal standard plane identification | VRA | 137 | CCA 4.0I |
| DDTI [57] | thyroid | Thyroid nodule identification Thyroid nodule localisation | VRA LL | 110 | - |
| The Open Kidney US Dataset [68] | kidney | Kidney detection Kidney Diag view identification | VRA OD | 110 | CC BY-NC-SA |
| FPUS23 [58] | Fetal abdomen Fetal arm Fetal head Fetal legs | Fetal diagnostic planes identification Fetal US report generation | VRA RP | 752 | MIT |
| Echogenic [19] | Fetal abdomen | Fetal abdominal organ detection | OD | 102 | CCA 4.0 |
| FALLMUD [24] | Crural muscles | Muscle detection | OD | 100 | - |
| Micro-US Prostate Segmentation Dataset [64] | Prostate | Prostate localisation Prostate Diag view identification | VRA LL | 110 | CCA 4.0I |
| CAMUS [39] | Heart ED Heart ES Heart 2CH Heart 4CH | Heart ejection fraction estimation Heart atrium and ventricle localisation | VRA OD CVE | 316 | CC BY-NC-SA 4.0 |
| Breast Lesion Detection in US Videos [44] | Breast benign Brest malignant | Breast lesion classification | Diag | 171 | - |
| Breast US Images Dataset [2] | Breast | Breast cancer level classification Breast tumour localisation Brest Diag view identification | Diag VRA LL | 210 | CC0: PD |
| Dermatologic Ultrasound Images for classification [38] | Skin | Skin tumor level classification | Diag | 100 | - |
| Polycystic Ovary Ultrasound Images Dataset [77] | Ovary | Polycystic Ovary Syndrome localisation | VRA | 10 | CC0: PDD |
| CUBS [50] | Carotid | Carotid thickness estimation Carotid detection Catotid Diag view identification | VRA OD CVE | 681 | CCA 4.0I |

Continued on next page

(Continued) Table 6

| Dataset | Anatomy | Clinical scenarios | Task | Case | License |
|---|---|--|--------------------|------|----------------|
| Knee US dataset in a population-based cohort [53] | Knee | Knee US KL and pain grad classification Knee Diag view identification Knee lesion localisation | Diag VRA OD | 326 | CC0 1.0 |
| HC18 [28] | Fetal head | Fetal head circumference estimation Fetal head detection | OD CVE | 202 | CCA 4.0I |
| KFGNet [52] | Thyroid | Thyroid nodule level classification Thyroid nodule localisation | Diag LL | 206 | - |
| Thyroid [36] | Thyroid Left Thyroid right | Thyroid Diag view identification | VRA | 563 | CC BY |
| GDPHSYSUCC [51] | Breast | Breast lesion classification | Diag | 109 | - |
| LEPset [43] | Pancreas | Pancreatic cancer classification | Diag | 101 | CCA 4.0I |
| COVID-BLUES [76] | Lung | COVID-19 level classification Lung US caption generation Lung Diag view identification | Diag VRA CG | 318 | ANN 4.0 I |
| Ultrasound Guided Regional Anesthesia [72] | Brachial plexus | Brachial plexus detection | OD | 179 | Non-commercial |
| Unity Imaging Collaborative [67] | Cardiac | Cardiac Keypoint Detection | KD | 500 | CCANN 4.0 I |
| C-TRUS Dataset [40] | Colon | Colon wall detection | OD | 166 | - |
| ACOUSLIC-AI [61] | Fetal abdominal | Fetal abdominal circumference estimation Fetal adominal OD | VRA OD CVE | 310 | CCANCSA 4.0I |
| PSFHS [7] | Fetal head Fetal pubic symphysis | Fetal head detection Fetal pubic symphysis detection | OD | 100 | CCA 4.0I |
| JNU-IFM [48] | Fetal head Fetal pubic symphysis | Fetal view identification Fetal head detection Fetal pubic symphysis detection | VRA OD | 202 | CC BY 4.0 |
| Dataset of B-mode fatty liver US images [13] | Liver | Liver steatosis classification Liver fat value estimation Liver Diag view identification | Diag VRA CVE | 222 | CCA 4.0I |
| African Fetal Standard Plane [63] | Fetal abdomen Fetal brain Fetal femur Fetal thorax | Fetal standard plane identification | VRA | 10 | CCA 4.0I |
| BrEaST [56] | Breast | Breast LL | LL | 100 | CC BY 4.0 |

Continued on next page

(Continued) Table 6

| Dataset | Anatomy | Clinical scenarios | Task | Case | License |
|--|---|---|-------------|------|------------|
| Ultrasound Images for Breast Cancer [60] | Breast | Breast cancer classification | Diag | 100 | CC0: PD |
| US simulation and segmentation [73] | Abdominal | Abdominal OD | OD | 100 | - |
| Carotid Artery Ultrasound and Color Doppler [55] | External carotid left carotid right carotid | Carotid Diag view identification | VRA | 100 | Apache 2.0 |
| AUITD [49] | Thyroid | Thyroid lesion classification | Diag | 100 | - |
| Auto-PCOS classification [49] | Ovary | Polycystic Ovary Syndrome classification Ploycystic Diag view identification | Diag VRA | 218 | CCA 4.0I |
| Auto-PCOS classification [49] | Ovary | Polycystic Ovary Syndrome classification | Diag | 100 | CC BY 4.0 |