

EchoVLM: Dynamic Mixture-of-Experts Vision-Language Model for Universal Ultrasound Intelligence

Chaoyin She¹, Ruifang Lu², Lida Chen^{2,*}, Wei Wang^{2,*}, Qinghua Huang^{1,3,*}

¹School of iOPEN, Northwestern Polytechnical University, Xi'an, China

²The First Affiliated Hospital of Sun Yat-Sen University, Guangzhou, China

³College of Mechanical Engineering, Tongji University, Shanghai, China

{chenlda, wangw73}@mail.sysu.edu.cn qinghua_huang@tongji.edu.cn qhhuang@nwpu.edu.cn

Abstract

Ultrasound is the preferred early cancer screening modality due to non-ionizing radiation, cost-effectiveness, and real-time imaging, yet conventional diagnosis relies heavily on physician expertise, causing significant subjectivity and limited efficiency. Vision-Language Models (VLMs) show promise but lack ultrasound-specific knowledge and multi-organ generalization. We propose EchoVLM, the first open-source 10-billion-parameter ultrasound-tailored VLM with a Mixture-of-Experts (MoE) architecture. It is infused with knowledge across seven anatomical systems, trained on 208,941 clinical cases, 1.47 million ultrasound key-frame images, and over 100 diseases or imaging findings. Supporting clinical report generation, diagnosis prediction, and Visual Question Answering (VQA), it outperforms Qwen2-VL by 7.58 BLEU-1 and 3.45 ROUGE-1 points in report generation. This work shows substantial potential for establishing a general-purpose ultrasound VLM and lays a technical foundation for clinical translation. Source code and model weights are available at <https://github.com/Asunatan/EchoVLM>.

1 Introduction

Ultrasound imaging has become a cornerstone of clinical diagnostics, distinguished by the absence of ionizing radiation, cost-effectiveness, and real-time dynamic visualization capabilities—attributes that render it indispensable for early cancer detection, prenatal care, and dynamic assessment of organ structure and function. However, conventional ultrasound diagnosis remains heavily dependent on radiologists' specialized expertise, with manual interpretation introducing inter-observer variability, diagnostic delays, and suboptimal treatment efficiency. While Visual Language Models (VLMs) have made significant advances in multimodal perception, their application to ultrasound diagnosis

*Corresponding authors

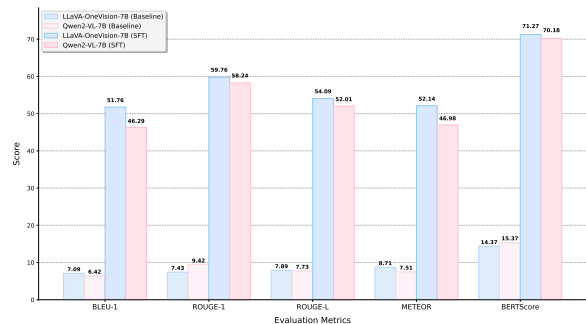


Figure 1: Comparative performance evaluation of generic and ultrasound-specialized VLMs.

faces a critical domain-specific limitation: general-purpose VLMs exhibit limited domain generalization to ultrasound medicine. Our preliminary experiments (Figure 1) reveal a substantial performance gap between baseline generic models and their supervised fine-tuned (SFT) variants, confirming that off-the-shelf VLMs cannot effectively capture the nuanced, domain-specific features of ultrasound images. This limitation directly restricts their clinical utility in ultrasound workflows, underscoring the urgent need to develop ultrasound-specialized VLMs.

To address this gap, we introduce EchoVLM, the first universal ultrasound-specialized VLM with tens of billions of parameters, founded on three core innovations: (1) We curated the largest multi-organ ultrasound dataset to date, covering seven anatomical systems based on 208,941 clinical cases from 15 hospitals, 1.47 million key-frame ultrasound images, and over 100 related diseases and imaging findings, ensuring comprehensive coverage for robust model training. (2) Inspired by Self-Instruct (Wang et al., 2023), we propose an expert-validated few-shot prompting mechanism to build a multi-task instruction-tuning data generation pipeline, which has generated 1.8 million pairs of instruction-tuning data by synthesizing diverse diagnostic scenarios while ensuring clinical

accuracy via expert oversight. (3) We utilize a Dual-path MoE module for knowledge injection while preserving pre-acquired knowledge; its dynamic routing mechanism enhances adaptability to ultrasound’s heterogeneous complexity, enabling task-specific subnetworks to specialize in distinct diagnostic domains for improved efficiency and precision. In summary, our key contributions are:

- We pioneer EchoVLM, the first universal ultrasound-specialized VLM with tens of billions of parameters tailored to address ultrasound diagnostic challenges.
- We curate a large-scale multi-organ ultrasound dataset (208,941 cases from 15 hospitals, 1.47 million key-frame images, over 100 diseases/imaging findings, covering seven anatomical systems) and develop an expert-validated multi-task instruction-tuning pipeline generating 1.8 million data pairs.
- We integrate a Dual-path MoE module into EchoVLM for knowledge injection while preserving pre-acquired knowledge, whose dynamic routing mechanism markedly enhances adaptability to ultrasound task complexity.

2 Related Work

Large language models (LLMs) have driven paradigm shifts in artificial intelligence, especially in natural language processing (NLP) for comprehension, generation, and text-based tasks. However, traditional unimodal LLMs inherently lack visual perception capabilities, which limits their applicability to real-world cross-modal scenarios. To address this, researchers first leveraged image-text alignment methods (Radford et al., 2021; Vasu et al., 2024; Zhang et al., 2024a) to establish a foundational bridge between visual and textual modalities; subsequently, semantic embedding layers (Liu et al., 2023; Lin et al., 2024a,b) and cross-attention layers (Li et al., 2023b; Alayrac et al., 2022) have been proposed to achieve effective visual-textual feature integration. Notably, the BLIP series pioneered cross-attention for dynamic visual-textual interaction, while LLaVA introduced visual instruction tuning to enhance perceptual capabilities via visual data fine-tuning. Subsequent advancements improved performance through dataset expansion (Liu et al., 2024a; Zhang et al., 2024d), higher image resolution (Liu et al., 2024b; Wang

et al., 2024b; Hong et al., 2024), multi-image/video understanding (Li et al., 2024a; Zhang et al., 2025), projection layer optimization (Li et al., 2026; Chen et al., 2025; Tong et al., 2024), and multi-visual encoders (Kar et al., 2024; Fan et al., 2024). Recent progress focuses on fine-grained understanding by integrating object localization (e.g., bounding boxes) and segmentation masks to boost spatial and semantic accuracy (Guo et al., 2024; Zhang et al., 2024b). Notably, VLMs have been increasingly applied in medicine, with LLaVA-Med (Li et al., 2023a), Medgemma (Sellergren et al., 2025), HuatuoGPT-Vision (Chen et al., 2024), and Lingshu (Team et al., 2025b) as representative examples. However, such systems face critical limitations in specialized clinical settings, especially in terms of insufficient ultrasound-specific knowledge that impairs structured report generation. Additionally, their limited context length fails to handle the multi-image scenarios common in ultrasound. To address these domain-specific issues, this study proposes an ultrasound-specific VLM architecture optimized for diagnostic workflows and standardized sonographic terminology.

3 Method

3.1 Data Collection and Instruction-Tuning Data Generation Pipeline

To develop a VLM tailored to ultrasound imaging, we compiled a comprehensive dataset from 15 hospitals, as shown in Figure 2. It covers seven major anatomical systems commonly assessed via ultrasound: liver, kidneys, thyroid, vascular system, gynecological organs, heart, and breasts. A rigorous data filtration protocol ensured quality: (1) Image Filtering: Only single-region images were extracted from hospital databases to avoid multi-region ambiguity. Images without corresponding reports were manually removed for alignment between imaging and clinical data. (2) Text Filtering: Sensitive patient information was deleted using regular expression matching to ensure data privacy, and irrelevant reports or those lacking imaging data were manually deleted. To advance ultrasound-specific VLMs, we pioneered the redefinition of data desensitization (see Appendix A for details). This process yielded 208,941 cases with 1.47 million key-frame images, covering over 100 diseases and imaging findings (Figure 3).

Leveraging this dataset, we established a structured instruction-tuning data generation pipeline

Dataset	Anatomical Coverage	Dataset Scale	Imaging Modality	Target Tasks	Multimodal
BUSI (Al-Dhabyani et al., 2020)	Breast	780 images 600 patients	2D B-mode	Segmentation Classification	✗
OASBUD (Piotrkowska-Wróblewska et al., 2017)	Breast	200 images 100 patients	RF Ultrasound 2D B-mode	Segmentation Classification QUS Analysis	✗
BUS-UCLM (Vallez et al., 2025)	Breast	683 images 38 patients	2D B-mode	Segmentation Classification	✗
TN3K (Gong et al., 2021)	Thyroid	3,493 images 2,421 patients	2D B-mode	Segmentation	✗
TNUI-2021 (Nie et al., 2022)	Thyroid	1,381 images 483 patients	2D B-mode	Segmentation Classification	✗
DDTI (Pedraza et al., 2015)	Thyroid	134 images 99 patients	2D B-mode	Segmentation Classification	✗
EchoNet-Dynamic (Ouyang et al., 2019)	Heart	10,030 videos 10,030 patients	Echocardiography Videos Color Doppler Spectral Doppler Tissue Doppler	Segmentation, Function Assessment	✗
EchoPrime (Vukadinovic et al., 2026)	Heart	12.1M videos 275K studies 67.8M text tokens	Echocardiography Videos Color Doppler Spectral Doppler Tissue Doppler	VLM Pre-training Classification Diagnosis Cross-modal Retrieval	✓
KMVE (Li et al., 2024b)	Breast Thyroid Liver	7,390 patients 7,390 reports	2D B-mode	VLM Pre-training Image Captioning Report Generation Description Generation	✓
FetalCLIP (Maani et al., 2025)	Obstetrics	210,035 images Paired clinical text	2D B-mode Color Doppler	VLM Pre-training Classification Segmentation Detection	✓
EchoCLIP (Christensen et al., 2024)	Heart	1.03M videos 99,870 patients 99,870 clinical report	Echocardiography Videos Color Doppler Spectral Doppler Tissue Doppler	VLM Pre-training Classification Regression Cross-modal Retrieval	✓
Sonomate (Guo et al., 2026)	Obstetrics	525 unique video and audio pairs 2.7M frames 63.8K sentences	Fetal Ultrasound Videos	VLM Pre-training Classification Detection Cross-modal Retrieval VQA	✓
EchoVLM (Ours)	7 Organ Systems	1.47M images 208K patients 1.8M instruction pairs	2D B-mode Color Doppler Spectral Doppler Tissue Doppler	Report Generation VQA Diagnosis	✓

Table 1: Summary of Medical Ultrasound Datasets.

using few-shot prompting (Figure 2). Medical experts developed 21 exemplary templates across diverse pathologies, each integrating three clinically simulative components: (1) Ultrasound reports: Detailed records of lesion location, size, morphology, echogenicity, and other image-derived features. (2) Ultrasound diagnosis: Standardized diagnostic summaries synthesizing key findings. (3) VQA pairs: Multidimensional sets covering interpretation, risk stratification, counseling, surveillance, and treatment planning. Templates were categorized into open- and closed-ended types with tailored prompting. For open-ended ones, models generated both questions and answers via example reference; for closed-ended ones, questions were model-generated and answers extracted from real reports to ensure validity. ROUGE-L and Simhash algorithms deduplicated data. The open-ended subset underwent dual accuracy validation: automated evaluation via a VLM/LLM pool and expert manual review of random samples. This pipeline produced 1.8 million high-quality instruction-tuning pairs (see Appendix A for details).

To further contextualize the contributions of our proposed dataset and instruction-tuning pipeline,

we conducted a rough statistical summary of existing studies in the ultrasound imaging domain, which are summarized in Table 1). This analysis demonstrates a clear paradigm shift in the field, wherein research is rapidly transitioning from traditional single-modality, task-specific approaches to large-scale, multimodal corpora that are tailored for vision-language pretraining and multi-task learning. Such an evolutionary trend further underscores the growing necessity of comprehensive, clinically aligned datasets to empower the development of next-generation ultrasound VLMs.

3.2 Architecture of EchoVLM

We introduce EchoVLM, a vision-language model tailored for clinical ultrasound analysis, developed via targeted domain specialization of the Qwen2-VL (Wang et al., 2024a) foundational model. Instead of naively fine-tuning the model components on ultrasound corpora, we introduce a Dual-path MoE mechanism which injects domain knowledge while avoiding destructive updates to pre-existing representations.

For the visual encoder, any RGB ultrasound frame $v \in \mathbb{R}^{H \times W \times 3}$ are encoded by a native dy-

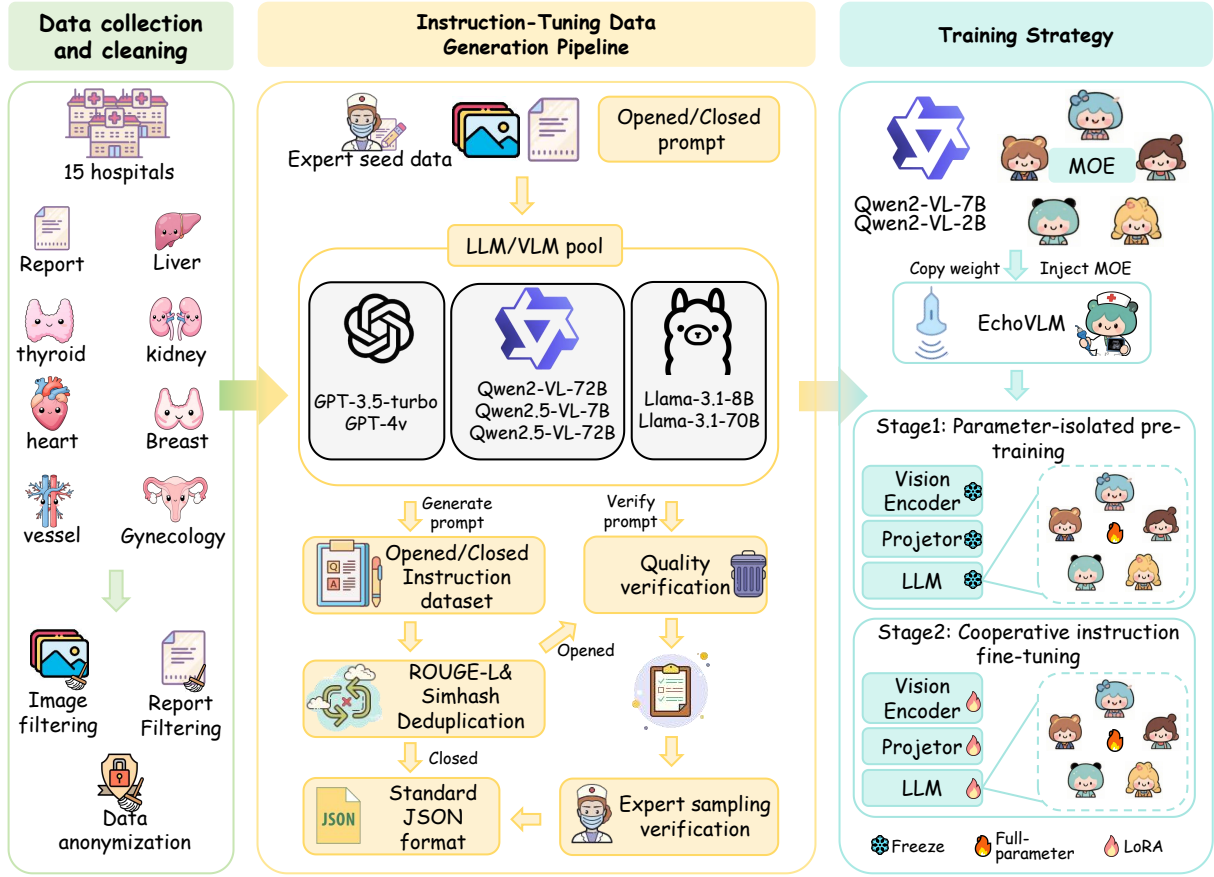


Figure 2: Overview of EchoVLM framework.

dynamic resolution Vision Transformer (ViT) to generate a discrete visual token representations $V = [v_1, v_2, \dots, v_m] \in \mathbb{R}^{M \times C}$. Here, $M = HW/14^2$ denotes the number of visual tokens derived from the original spatial resolution $H \times W$. These visual tokens are subsequently transformed into the dimensional space of the Qwen-2 via a Multi-layer Perceptron (MLP) projector $f(\cdot)$, yielding dimensionally-mapped visual embeddings $V \in \mathbb{R}^{M \times D}$. Simultaneously, textual inputs undergo parallel processing: the input prompts are first tokenized and then embedded through the word-embedding layer $w(\cdot)$, generating the textual token sequence $T = [t_1, t_2, \dots, t_N] \in \mathbb{R}^{N \times D}$, where N represents the number of text tokens determined by the input prompt length. The concatenated representation $X_0 = [V; T] \in \mathbb{R}^{(M+N) \times D}$ is subsequently fed into the LLM which comprises multiple Transformer blocks. Each block integrates multi-head self-attention (MSA), RMS normalization (RMSNorm), residual connections, and the Dual-path MoE block, with its processing workflow formulated as follows:

$$X'_i = MSA(RMSNorm(X_{i-1})) + X_{i-1} \quad (1)$$

$$X_i = MoE(RMSNorm(X'_i)) + X'_i \quad (2)$$

3.3 Dual-path MoE

Structurally, the Dual-path MoE layer comprises two complementary sets of experts. First, a static expert is instantiated by copying the original Qwen2 Feed-Forward Network (FFN) and immediately frozen; its parameters therefore act as a resilient anchor that conserves generic semantic capacity. Second, a battery of active experts is appended and trained. Within this group we further distinguish (1) a shared expert (S) that processes every token, thereby sustaining a universal ultrasound representation, and (2) a cohort of routing experts (E) that are sparsely activated via a top-2 gating function conditioned on token-level features.

$$Y = \alpha FFN(X) + (1 - \alpha) \left[\lambda S(X) + \sum_{i=1}^k g_i(X) E_i(TopK(X)) \right] \quad (3)$$

$$g_i(X) = \frac{e^{f_i(X)}}{\sum_{j=1}^k e^{f_j(X)}} \quad (4)$$

In these equations, g_i quantifies the contribution of expert E_i , f_i denotes the routing logits, and α is a learnable scalar that modulates the equilibrium be-

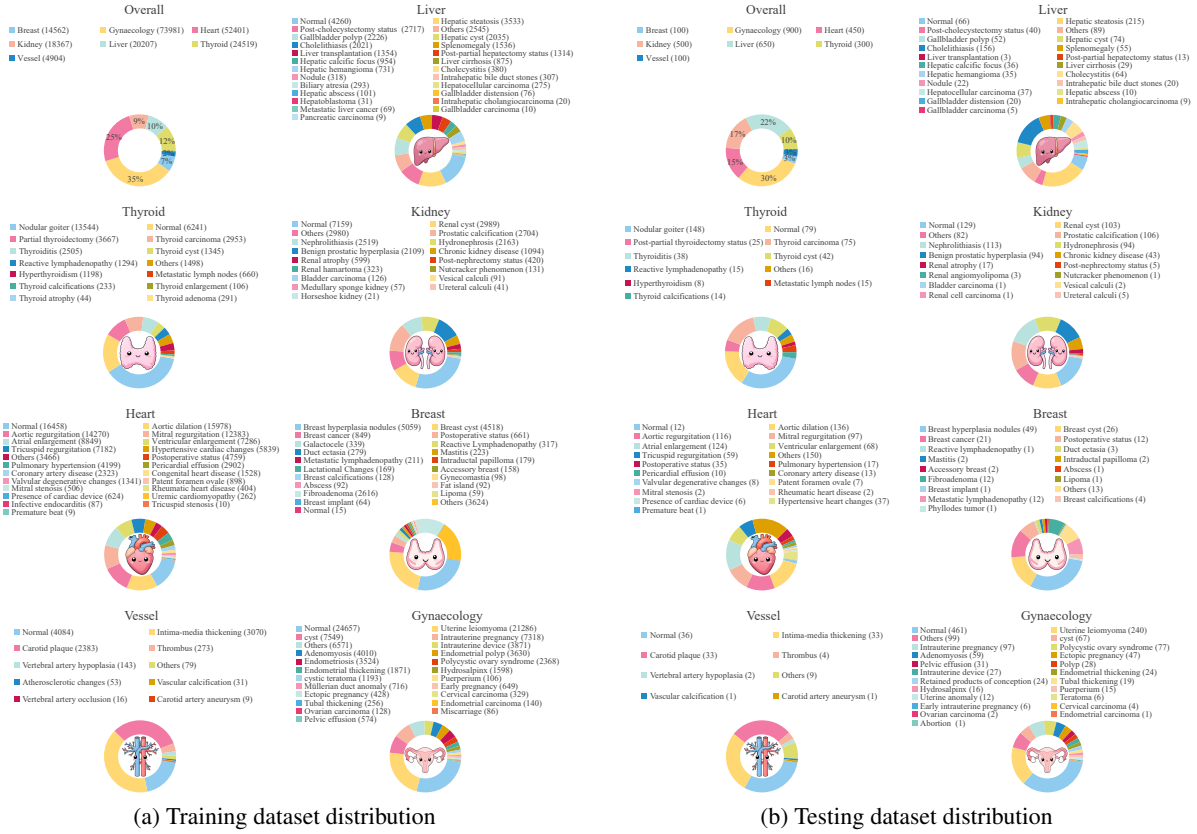


Figure 3: Distribution of cases across seven anatomical systems in a multicenter dataset, with numbers in parentheses indicating case counts. Note that a single case may involve multiple diseases.

tween generic knowledge and ultrasound-specific information.

3.4 Training Strategy

3.4.1 Stage I: Parameter-isolated pre-training.

All original parameters of the Qwen2-VL are completely frozen, and only the newly-introduced MoE blocks are activated. This isolation strategy prevents catastrophic forgetting of generic multimodal knowledge, steering the expert layers to exclusively acquire domain-specific ultrasound representations.

3.4.2 Stage II: Cooperative instruction fine-tuning.

We introduce a collaborative optimization mechanism by applying low-rank adaptation (LoRA: $\mathbf{W} = \mathbf{W}_0 + \Delta\mathbf{W}$, where $\Delta\mathbf{W} = \mathbf{A}\mathbf{B}^\top$ with $\mathbf{A} \in \mathbb{R}^{d \times r}$, $\mathbf{B} \in \mathbb{R}^{d \times r}$) (Hu et al., 2021) for lightweight parameter tuning of the base model, while maintaining full parameter updates for the MoE components. Notably, for the architecture comprising frozen FFN layers and activated MoE modules, we develop dynamic modulation

parameters $\lambda \in [0, 1]$ to balance the contributions between general world knowledge and ultrasound domain knowledge.

3.4.3 Training Objectives.

We optimize the model by minimizing a compound loss that balances the auto-regressive loss with an expert-load-balancing penalty. The vision encoder processes input images to generate m visual tokens V . For textual inputs, raw text undergoes tokenization and is subsequently mapped to dense embeddings using an embedding layer, yielding n textual tokens T . Their concatenation forms the model input $X = [V, T]$ of length $K = m + n$. During training, the model autoregressively generates a response of length L ; the autoregressive loss is thus defined as:

$$\mathcal{L}_{ar} = - \sum_{i=1}^L \log p_{\theta}(y_i | X, y_{<i}) \quad (5)$$

where θ is a trainable parameter. To ensure balanced utilization of the E experts, we introduce an auxiliary load-balancing loss. Let $g_e(t) \in [0, 1]$ denote the gating weight for expert e and token t ,

where $\sum_{e=1}^E g_e(t) = 1$. For a mini-batch of tokens \mathcal{B} , we compute the dispatch ratio F_e , representing the fraction of tokens routed to expert e :

$$F_e = \frac{1}{|\mathcal{B}|} \sum_{t \in \mathcal{B}} \mathbb{I}[g(t) \rightarrow e] \quad (6)$$

where $g(t) \rightarrow e$ denotes that token t is routed to expert e , $\mathbb{I}[\cdot]$ is the indicator function that equals 1 if token t is routed to expert e , and 0 otherwise. The average gating probability is defined as:

$$G_e = \frac{1}{|\mathcal{B}|} \sum_{t \in \mathcal{B}} g_e(t). \quad (7)$$

The auxiliary balancing loss is then given by:

$$\mathcal{L}_{bal} = \sum_{e=1}^E F_e G_e. \quad (8)$$

The overall training objective combines both components:

$$\mathcal{L}_{total} = \mathcal{L}_{ar} + \gamma \mathcal{L}_{bal} \quad (9)$$

where γ is a hyperparameter controlling the strength of the load-balancing penalty.

4 Experiments

4.1 Implementation Details

Following LLaVA-Med’s protocol, the two-stage training framework adopts phase-specific datasets. Stage I initializes newly introduced MoE modules using 208,941 clinical reports and 1.47 million ultrasound key frames, exclusively training them to capture domain-specific visual-textual patterns without modifying the base model’s pre-trained parameters. Stage II performs Cooperative Instruction Fine-tuning on 1.8 million instruction-following samples, integrating the base model and MoE modules to enhance instruction understanding and response capabilities. For reproducibility, the model is evaluated on a held-out test set (27,577 ultrasound images and 3,000 reports), with greedy decoding during inference ensuring deterministic and reproducible outputs. More details are provided in the Appendix B.

4.2 Results

4.2.1 Report Generation

Table 2 compares 12 VLMs (general and medical-specific) across seven anatomical systems using five key metrics. EchoVLM (11B) is the core highlight, consistently outperforming counterparts (e.g.,

Qwen2-VL (Wang et al., 2024a), Qwen2.5-VL (Bai et al., 2025), Gemma-3 (Team et al., 2025a)) and achieving SOTA in multiple scenarios. It tops BERTScore in four systems (breast:73.00, kidney:75.89, etc.), leads ROUGE-1/ROUGE-L in breast, kidney, liver, and thyroid, and ranks first in BLEU-1 and METEOR for most domains. Its average scores are outstanding, with the highest scores in BLEU-1 (53.87), ROUGE-1 (61.69), and ROUGE-L (55.78), and the second-highest in METEOR and BERTScore. This proves its strong domain-specific performance and generalizability. Even the smaller EchoVLM (3B) performs admirably (e.g., 3rd in breast BLEU-1/ROUGE-1). In contrast, compared with other models adopting full-parameter fine-tuning, this outstanding performance of EchoVLM demonstrates the high efficiency of leveraging MoE to inject domain knowledge.

Traditional n-gram overlap and text similarity metrics (BLEU, ROUGE, BERTScore) are widely used for report generation evaluation but correlate imperfectly with clinical correctness, as they favor surface-level text matching over accurate medical concepts, attributes and critical findings. Therefore, we supplement our experiments with fine-grained entity-level metrics, extracting anatomical structures, findings, attributes and statuses via a clinically validated extractor and computing four standard classification metrics. Complete results are provided in Appendix F.

4.2.2 Ultrasound Diagnosis

Ultrasound diagnosis is a highly specialized synthesis of sonographic findings, requiring precise integration of anatomical and imaging findings with high professional accuracy. To validate EchoVLM’s efficacy in this task, we evaluated it against other models across seven anatomical systems. Table 3 shows EchoVLM (11B) outperforms other general and medical-specific models in most categories, achieving SOTA in key systems like kidney (BLEU-1: 77.56, BERTScore: 87.03) and liver (ROUGE-1: 79.87, BERTScore: 82.85). It also leads all average metrics (BERTScore:75.44, ROUGE-1:72.51), reflecting strong alignment with clinical terminology and superior diagnostic reasoning.

4.2.3 Vision Question Answering

Ultrasound VQA, as an open-ended question task, typically demands models with hierarchical com-

Anatomical	Metric↑	Medgemma (4B)	LLaVA1.5 (7B)	HuatuogPT-Vision (7B)	Lingshu (7B)	LLaVA-OneVision (7B)	Qwen2-VL (7B)	LLaVA-Med (7B)	Qwen2.5-VL (7B)	LLaVa-NeXT (13B)	Gemma-3 (12B)	EchoVLM (3B)	EchoVLM (11B)
Breast	BLEU-1	21.91	46.49	50.10	48.99	49.59	41.74	51.17	50.55	49.50	48.63	50.55	50.76
	ROUGE-1	41.13	57.07	60.97	60.87	61.45	62.85	60.44	61.78	59.03	60.49	62.24	64.38
	ROUGE-L	34.11	49.33	51.64	52.18	51.47	53.46	52.03	52.38	50.53	49.59	53.31	55.21
	METEOR	26.42	45.29	49.76	49.32	49.86	47.32	48.41	49.51	47.11	47.44	49.8	50.14
	BERTScore	51.36	66.92	70.35	70.39	70.84	72.21	69.50	71.53	68.32	70.35	70.82	73.00
Gynecology	BLEU-1	46.33	43.51	48.75	51.40	47.83	49.88	42.63	47.28	43.01	49.05	45.79	52.52
	ROUGE-1	55.68	48.80	53.18	56.64	54.98	56.97	51.64	52.15	52.06	55.71	53.15	59.19
	ROUGE-L	48.27	42.58	44.19	48.33	47.52	47.92	45.12	44.01	45.59	47.50	46.29	51.64
	METEOR	44.63	42.97	47.55	49.96	47.88	49.36	44.32	45.52	44.80	48.32	45.94	52.76
	BERTScore	65.11	61.19	63.95	66.54	66.01	65.93	62.02	63.34	62.28	65.51	63.14	67.65
Heart	BLEU-1	71.76	57.06	73.52	77.45	77.43	72.25	70.04	77.72	72.24	77.97	72.18	76.48
	ROUGE-1	74.74	64.65	74.17	75.94	73.49	75.06	71.71	75.09	73.81	76.82	75.94	78.18
	ROUGE-L	77.12	65.50	77.04	79.92	78.35	77.56	74.99	79.54	76.76	80.35	77.84	80.64
	METEOR	63.13	50.86	66.32	73.30	72.07	63.68	61.77	73.02	64.67	73.78	63.48	70.73
	BERTScore	87.64	82.52	87.82	89.94	89.40	87.90	85.27	89.86	86.20	90.12	87.91	89.80
Kidney	BLEU-1	45.72	39.79	48.51	52.09	49.54	43.80	41.71	44.84	42.27	49.70	43.79	59.23
	ROUGE-1	57.69	46.28	54.77	58.32	59.49	49.23	51.86	48.23	52.59	55.75	50.56	65.78
	ROUGE-L	47.66	38.36	43.73	49.26	50.82	41.95	43.53	40.38	44.25	47.47	42.07	57.14
	METEOR	44.15	39.68	47.70	52.57	51.57	42.82	42.81	43.03	43.38	49.87	43.17	59.58
	BERTScore	67.96	58.86	66.86	70.63	70.23	62.65	63.37	63.08	63.91	68.31	63.20	75.89
Liver	BLEU-1	47.00	40.19	49.40	53.11	53.87	40.48	43.10	44.28	43.68	53.64	50.91	58.01
	ROUGE-1	61.90	44.34	55.02	60.95	60.59	51.93	53.77	47.41	54.53	59.82	63.57	67.06
	ROUGE-L	58.28	40.18	50.85	56.67	56.35	46.78	50.89	43.39	51.68	53.89	59.30	62.40
	METEOR	48.93	38.85	47.82	52.87	53.43	41.54	45.63	41.52	46.39	53.01	53.59	58.44
	BERTScore	70.86	58.82	67.73	71.58	72.22	65.98	64.66	62.33	65.22	70.91	72.51	75.68
Vessel	BLEU-1	39.10	29.75	41.15	40.65	42.54	37.62	16.12	37.39	29.56	42.53	31.64	29.54
	ROUGE-1	65.14	45.68	53.40	63.78	60.81	60.63	20.06	55.57	45.65	63.51	42.14	38.15
	ROUGE-L	58.36	39.15	46.56	56.98	53.79	55.38	13.36	48.43	39.06	56.86	35.95	32.29
	METEOR	50.80	36.98	47.39	50.80	50.77	46.50	15.27	45.06	36.97	51.28	34.41	30.38
	BERTScore	74.00	54.97	62.27	72.88	69.03	70.01	28.23	64.32	55.08	71.74	51.76	48.07
Thyroid	BLEU-1	23.17	26.79	39.68	45.59	41.50	38.25	44.27	44.23	44.00	39.37	45.13	50.55
	ROUGE-1	38.35	35.21	48.72	54.60	47.51	51.03	53.30	52.93	52.92	48.91	52.65	59.10
	ROUGE-L	30.50	27.81	40.54	45.50	40.30	41.00	44.45	44.10	43.99	39.68	44.47	51.16
	METEOR	25.52	26.57	40.99	44.85	39.41	37.63	45.51	43.95	45.13	38.13	44.21	49.87
	BERTScore	50.13	49.51	59.75	66.06	61.15	66.59	63.86	64.89	63.54	61.36	64.11	69.54
Average	BLEU-1	42.14	40.51	50.16	52.75	51.76	46.29	44.15	49.47	46.32	51.56	48.57	53.87
	ROUGE-1	56.38	48.86	57.18	61.59	59.76	58.24	51.83	56.17	55.80	60.14	57.18	61.69
	ROUGE-L	50.61	43.27	50.65	55.55	54.09	52.01	46.34	50.32	50.27	53.62	51.32	55.78
	METEOR	43.37	40.17	49.65	53.38	52.14	46.98	43.39	48.80	46.92	51.69	47.80	53.16
	BERTScore	66.72	61.83	68.39	72.57	71.27	70.18	62.42	68.48	66.36	71.19	67.64	71.38

Table 2: Comparison results for report generation.

Anatomical	Metric↑	Medgemma (4B)	LLaVA1.5 (7B)	HuatuogPT-Vision (7B)	Lingshu (7B)	LLaVA-OneVision (7B)	Qwen2-VL (7B)	LLaVA-Med (7B)	Qwen2.5-VL (7B)	LLaVa-NeXT (13B)	Gemma-3 (12B)	EchoVLM (3B)	EchoVLM (11B)
Breast	BLEU-1	61.94	54.79	65.29	70.37	65.46	67.63	69.84	68.09	68.82	67.91	67.92	71.36
	ROUGE-1	77.42	70.42	77.00	78.74	75.78	78.13	78.37	77.61	77.61	75.19	78.11	80.77
	ROUGE-L	70.17	58.71	69.93	72.98	67.22	71.38	72.30	69.46	71.21	68.28	70.00	76.04
	METEOR	65.99	64.29	71.25	75.30	73.99	73.88	74.34	74.92	73.53	72.79	78.78	74.99
	BERTScore	77.81	69.20	77.07	79.23	75.70	78.32	78.72	77.34	77.98	76.70	78.03	81.68
Gynecology	BLEU-1	28.38	29.77	39.59	36.46	33.46	40.46	37.80	36.37	38.21	37.42	43.32	48.15
	ROUGE-1	47.43	38.66	48.71	45.70	45.33	53.13	53.56	46.39	53.83	48.87	56.24	62.82
	ROUGE-L	42.44	36.29	43.74	42.01	40.79	48.44	49.44	42.25	49.74	44.99	52.15	58.25
	METEOR	33.43	36.92	44.62	42.28	40.16	47.32	45.82	43.53	46.19	44.12	50.86	55.95
	BERTScore	52.65	43.56	53.50	50.97	50.18	57.81	58.53	51.38	58.82	53.99	60.03	66.58
Heart	BLEU-1	55.99	42.41	55.32	64.60	61.00	66.26	67.77	60.12	69.31	65.04	66.30	69.62
	ROUGE-1	71.58	50.32	62.77	76.49	70.95	77.16	79.13	71.90	79.94	75.89	78.74	81.33
	ROUGE-L	64.63	44.81	55.69	69.31	63.39	69.78	71.26	64.78	72.55	68.57	71.47	74.35
	METEOR	59.00	43.54	59.43	67.68	64.16	68.59	70.26	63.46	71.63	68.02	68.91	72.63
	BERTScore	73.15	55.58	67.65	76.34	71.59	77.63	78.86	72.96	79.88	76.02	78.06	79.94
Kidney	BLEU-1	66.68	44.09	62.52	69.15	71.76	68.71	70.76	63.05	70.74	69.95	74.12	77.56
	ROUGE-1	77.06	52.81	70.77	76.37	79.32	76.87	78.10	71.77	78.00	78.01	81.70	83.42
	ROUGE-L	67.55	44.51	59.51	67.23	70.77	67.23	68.58	61.64	68.59	68.02	72.89	74.86
	METEOR	67.15	52.09	68.75	73.26	75.19	72.57	74.80	69.18	74.75	74.75	76.56	80.90
	BERTScore	81.51	60.40	76.07	80.80	83.46	81.46	82.17	77.08	82.13	81.78	84.71	87.03
Liver	BLEU-1	60.39	50.19	61.98	67.11	66.60	63.23	62.42	60.77	62.93	58.88	68.45	74.23
	ROUGE-1	70.87	59.10	69.42	73.84	73.02	73.63	70.45	68.58	70.80	65.89	75.85	79.87
	ROUGE-L	64.49	52.72	63.56	69.00	68.26	66.49	64.65	62.91	65.14	58.93	69.32	75.06
	METEOR	63.96	58.21	69.15	72.34	71.22	68.75	71.02	70.02	71.48	68.71	73.77	78.49
	BERTScore	73.55	63.36	73.24	77.40	77.01	76.14	73.65	73.06	73.96	70.62	77.81	82.85
Vessel	BLEU-1	53.41	31.68	49.23	47.86	49.65	59.56	27.54	42.85	35.88	46.86	38.70	36.27
	ROUGE-1	66.86	47.47	57.85	60.58	60.12	70.69	42.75	56.73	51.35	58.23	51.00	47.63
	ROUGE-L	65.32	44.99	55.49	59.09	58.31	69.54	36.99	54.51	49.05	56.41	46.33	44.08
	METEOR	59.53	40.97	56.04	53.12	55.46	64.98	45.73	53.06	44.78	51.32	45.06	40.74
	BERTScore	74.45	58.16	66.90	69.33	67.44	76.82	46.97	65.74	61.19	67.10	58.43	53.86
Thyroid	BLEU-1	49.51	32.60	54.82	52.93	53.03	51.86	52.27	55.24	52.17	52.38	57.23	62.49
	ROUGE-1	59.95	42.84	65.57	63.85	62.01	65.01	61.33	63.32	60.89	59.82	66.30	71.74
	ROUGE-L	54.39	36.84	61.37	58.85	56.08	59.56	57.12	58.82	56.78	53.76	60.97	67.46
	METEOR	58.29	44.49	63.79	62.23	53.03	61.06	61.90	63.65	61.47	61.17	64.29	70.03
	BERTScore	65.05	48.68	70.74	68.90	66.89	69.69	67.18	68.36	66.81	65.12	70.49	76.11
Average	BLEU-1	53.76	40.79	55.54	58.35	57.28	59.67	55.49	55.21	56.87	56.92	59.43	62.81
	ROUGE-1	67.31	51.66</										

Anatomical	Metric↑	Medgemma (4B)	LLaVA1.5 (7B)	HuatuogPT-Vision (7B)	Lingshu (7B)	LLaVA-OneVision (7B)	Qwen2-VL (7B)	LLaVA-Med (7B)	Qwen2.5-VL (7B)	LLaVA-NeXT (13B)	Gemma-3 (12B)	EchoVLM (3B)	EchoVLM(11B)
Breast	BLEU-1	25.03	25.81	32.35	36.43	36.79	33.60	25.61	36.62	28.31	35.86	34.09	35.80
	ROUGE-L	35.36	38.82	38.25	44.26	43.31	42.61	36.54	44.75	36.97	43.16	43.82	44.55
	ROUGE-L	27.18	30.72	29.82	36.13	35.03	34.73	26.45	36.75	28.48	34.57	35.62	36.31
	METEOR	23.73	23.49	29.37	30.77	31.47	28.19	22.00	31.07	25.59	29.26	29.09	29.76
	BERTScore	42.22	46.86	46.26	54.37	54.33	52.73	47.50	55.25	44.34	53.45	53.67	54.46
Gynecology	BLEU-1	23.33	21.58	30.17	32.76	29.88	28.02	24.46	31.64	26.23	33.12	25.74	25.75
	ROUGE-1	38.64	32.93	38.24	41.22	39.01	40.40	38.03	40.12	39.95	41.25	39.89	42.10
	ROUGE-L	30.01	24.20	28.62	31.64	30.01	31.69	28.87	31.24	31.50	31.61	31.09	33.30
	METEOR	22.52	20.29	25.99	29.20	27.13	25.21	22.76	28.07	24.01	30.04	24.2	25.59
	BERTScore	48.24	41.12	47.73	50.86	48.81	49.84	49.35	49.77	49.49	51.27	49.13	52.21
Heart	BLEU-1	21.46	21.58	27.77	27.98	30.43	22.96	21.20	28.18	23.75	28.88	20.10	23.55
	ROUGE-1	34.39	30.23	35.14	37.00	35.50	34.92	34.04	36.43	35.71	36.95	35.60	35.93
	ROUGE-L	25.88	22.56	26.54	28.87	27.69	26.75	24.92	28.51	27.53	28.84	27.28	27.62
	METEOR	21.84	20.66	24.85	25.72	27.76	22.87	20.22	25.73	23.06	26.43	22.10	23.86
	BERTScore	43.79	38.73	45.21	46.92	46.25	45.37	44.78	46.07	45.50	47.23	45.47	46.38
Kidney	BLEU-1	25.39	24.27	29.62	33.49	31.30	28.99	24.16	32.08	27.01	33.10	26.85	31.59
	ROUGE-1	39.60	35.74	39.16	42.75	41.46	41.23	37.07	41.51	40.45	42.57	40.98	42.47
	ROUGE-L	30.69	27.00	29.89	33.61	32.46	32.25	27.14	32.46	31.27	33.31	31.76	33.59
	METEOR	24.16	23.23	25.91	30.30	28.58	26.65	22.39	29.38	25.30	30.56	25.30	28.48
	BERTScore	48.79	44.78	49.09	52.93	51.67	51.10	48.49	51.67	49.86	52.93	49.91	52.24
Liver	BLEU-1	29.24	24.82	35.95	35.70	35.30	32.00	27.84	34.21	30.45	33.18	31.43	37.26
	ROUGE-1	39.75	34.46	42.16	42.51	40.65	41.24	37.73	40.86	40.13	39.55	41.59	44.36
	ROUGE-L	29.72	25.39	32.03	32.66	30.66	31.67	27.20	31.44	30.26	29.09	31.96	34.37
	METEOR	25.61	22.26	31.26	30.50	30.96	27.16	23.34	29.10	25.94	27.79	27.05	31.17
	BERTScore	49.05	43.10	52.44	52.59	51.34	51.19	48.94	51.13	49.45	49.39	51.02	54.20
Vessel	BLEU-1	28.83	20.97	32.10	35.24	34.27	29.71	24.30	32.99	21.48	34.25	26.46	27.06
	ROUGE-1	40.24	33.51	38.74	42.17	39.17	40.70	35.38	40.56	41.32	41.32	36.75	37.82
	ROUGE-L	30.82	24.89	29.63	33.56	31.30	32.20	25.30	31.87	25.30	32.61	27.72	28.79
	METEOR	25.85	20.77	26.52	29.40	28.94	25.63	22.02	28.65	20.94	29.14	24.18	24.55
	BERTScore	51.25	44.56	50.54	53.65	51.50	51.86	47.33	51.92	44.88	53.06	47.87	49.42
Thyroid	BLEU-1	25.14	26.30	31.93	33.41	34.77	30.69	25.07	33.65	25.81	31.73	28.3	29.92
	ROUGE-1	35.98	36.54	38.84	42.17	41.67	41.78	36.20	41.72	35.43	40.39	40.97	42.31
	ROUGE-L	28.36	29.17	30.99	34.44	34.23	34.25	27.35	33.89	27.61	32.45	33.34	34.74
	METEOR	24.31	23.74	28.95	29.74	30.76	27.38	22.28	29.85	24.48	27.39	26.43	27.86
	BERTScore	44.54	47.30	48.35	53.29	53.74	52.73	47.90	53.07	44.14	51.50	51.80	53.38
Average	BLEU-1	25.49	23.62	31.41	33.57	33.25	29.42	24.66	32.77	26.15	32.87	27.57	30.13
	ROUGE-1	37.71	34.60	38.65	41.73	40.11	40.41	36.43	40.85	37.50	40.74	39.94	41.36
	ROUGE-L	28.95	26.28	29.65	32.99	31.63	31.93	26.75	32.31	28.85	31.78	31.25	32.67
	METEOR	24.00	22.06	27.55	29.38	29.37	26.16	22.14	28.84	24.19	28.66	25.48	27.32
	BERTScore	46.84	43.78	48.52	52.09	51.09	50.69	47.76	51.27	46.81	51.26	49.84	51.76

Table 4: Comparison results for VQA.

tion and ultrasound diagnosis, with BLEU-1 gains of +4.58 and +3.48, ROUGE-L increases of +3.45 and +5.49, and BERTScore improvements of +1.27 and +4.06, respectively. These results indicate that shared experts enhance cross-task knowledge transfer and multimodal coherence in tasks requiring complex semantic synthesis.

Task	Metric↑	w/o Share Expert	w Share Expert
Report Generation	BLEU-1	49.29	53.87 (+4.58)
	ROUGE-1	58.42	61.69 (+3.27)
	ROUGE-L	52.33	55.78 (+3.45)
	METEOR	49.54	53.16 (+3.62)
	BERTScore	70.11	71.38 (+1.27)
Ultrasound Diagnosis	BLEU-1	59.33	62.81 (+3.48)
	ROUGE-1	67.46	72.51 (+5.05)
	ROUGE-L	61.67	67.16 (+5.49)
	METEOR	64.92	67.68 (+2.76)
	BERTScore	71.38	75.44 (+4.06)
VQA	BLEU-1	33.15	30.13 (-3.02)
	ROUGE-1	38.61	41.36 (+2.75)
	ROUGE-L	30.49	32.67 (+2.18)
	METEOR	29.09	27.32 (-1.77)
	BERTScore	49.60	51.76 (+2.16)

Table 5: Ablation study of the shared expert.

4.3.2 Impact of the Top-K Routing.

To investigate the impact of expert activation mechanisms within the MoE architecture, we implement two widely adopted routing strategies: Top-1 and Top-2 activation. As shown in Table 6, ablation studies demonstrate that Top-2 routing generally achieves superior performance in complex multimodal tasks compared to Top-1 routing, although

the magnitude of improvement varies across task characteristics. For report generation tasks, Top-2 routing enhances all evaluation metrics (e.g., +3.94 in BLEU-1, +2.81 in ROUGE-L, and +3.29 in METEOR). For ultrasound diagnostic tasks, it delivers remarkable gains across key metrics, including +4.10 in BLEU-1, +4.76 in ROUGE-1, and a substantial +5.00 in ROUGE-L. These results indicate that activating two specialized experts facilitates more effective knowledge integration and contextual understanding.

Task(4 Experts)	Metric↑	Top1	Top2
Report Generation	BLEU-1	49.93	53.87 (+3.94)
	ROUGE-1	58.80	61.69 (+2.89)
	ROUGE-L	52.97	55.78 (+2.81)
	METEOR	49.87	53.16 (+3.29)
	BERTScore	70.54	71.38 (+0.84)
Ultrasound Diagnosis	BLEU-1	58.71	62.81 (+4.10)
	ROUGE-1	67.75	72.51 (+4.76)
	ROUGE-L	62.16	67.16 (+5.00)
	METEOR	64.73	67.68 (+2.94)
	BERTScore	72.20	75.44 (+3.23)
VQA	BLEU-1	33.58	30.13 (-3.45)
	ROUGE-1	38.81	41.36 (+2.55)
	ROUGE-L	30.68	32.67 (+1.99)
	METEOR	29.60	27.32 (-2.28)
	BERTScore	49.96	51.76 (+1.80)

Table 6: Ablation study of Top-K routing strategies.

4.3.3 Impact of the Number of Routing Experts.

The number of experts represents a critical hyperparameter in EchoVLM where prior studies have demonstrated that scaling this parameter yields per-

Task	Metric \uparrow	0 Expert	2 Experts	4 Experts	6 Experts
Report Generation	BLEU-1	43.72	50.33	53.87	54.03
	ROUGE-1	56.92	58.75	61.69	61.91
	ROUGE-L	49.77	52.74	55.78	56.15
	METEOR	44.76	50.33	53.16	53.59
	BERTScore	70.17	70.51	71.38	71.71
Ultrasound Diagnosis	BLEU-1	58.70	59.66	62.81	63.23
	ROUGE-1	70.52	68.02	72.51	73.11
	ROUGE-L	64.52	62.45	67.16	67.94
	METEOR	64.05	65.59	67.68	68.47
	BERTScore	73.99	72.16	75.44	75.49
VQA	BLEU-1	27.64	33.52	30.13	33.73
	ROUGE-1	37.95	38.88	41.36	41.86
	ROUGE-L	28.96	30.69	32.67	32.68
	METEOR	26.53	29.46	27.32	29.64
	BERTScore	48.77	49.99	51.76	52.04
Training time (h/epoch)		46.99	52.12	55.36	59.10

Table 7: Ablation study of routing expert numbers

formance benefits. However its applicability within the ultrasound domain remains underexplored necessitating systematic investigation through ablation studies. As shown in Table 7, to address this gap we conducted an ablation experiment expanding the expert count from 0 to 6 specifically to determine whether the capacity gains observed in general VLMs could be effectively transferred to ultrasound-specific tasks characterized by distinct modality interactions and diagnostic complexity. Increasing from 0 to 4 experts delivers notable gains in high-complexity tasks report generation 10.15 BLEU-1 improvement ultrasound diagnosis 4.11 BLEU-1 improvement via improved model capacity enabling fine-grained cross-modal alignment. Compared to the 4-expert configuration six experts only yield marginal gains e.g. report generation BLEU-1 53.87 to 54.03 ultrasound diagnosis BLEU-1 62.81 to 63.23 but drastically increase training time 4 experts 55.36 h/epoch 6 experts 59.10 h/epoch. In contrast the 4-expert configuration achieves an optimal trade-off balancing superior performance and feasible training costs.

5 Conclusion

This study introduces EchoVLM, the first open-source 10B-parameter universal ultrasound-specialized VLM, with three core contributions: (1) curation of the largest multi-organ ultrasound dataset (208,941 clinical cases, 1.47M images across 7 anatomical systems); (2) development of a novel expert-validated few-shot prompting mechanism enabling a robust multi-task instruction fine-tuning pipeline; (3) integration of a Dual-path MoE architecture with dynamic routing, significantly enhancing adaptability to ultrasound’s heterogeneous imaging features. Experiments confirm EchoVLM outperforms SOTA across multiple anatomical systems. Limitations include dataset long-tail distribu-

tion and complex visual question answering requiring multi-step reasoning. Future work will focus on data rebalancing, expanded anatomical coverage, longitudinal patient data integration, and routing refinement to advance precise, efficient AI-assisted ultrasound diagnostics.

6 Limitations

While EchoVLM achieves promising overall performance, it exhibits suboptimal performance in vascular analysis. We hypothesize that this limitation stems from the long-tailed distribution of the entire dataset, where vascular cases account for the smallest proportion. This data scarcity likely led domain experts to prioritize dominant anatomical patterns from majority classes (e.g., breast/liver), inadvertently marginalizing subtle vascular features. Collectively, these findings demonstrate EchoVLM’s robustness in analyzing prevalent anatomical structures, while underscoring that its performance in vascular analysis could be significantly enhanced through data rebalancing strategies. Additionally, in open-ended tasks such as VQA, EchoVLM underperforms Lingshu (7B) in certain metrics, indicating potential room for improvement in its ability to handle the unstructured, open-ended nature of VQA tasks compared to general medical VLMs.

7 Acknowledgments

This work was supported by the National Natural Science Foundation of China (NSFC) Joint Fund for Regional Innovation and Development (Grant No. U25A20448), the National Natural Science Foundation of China (NSFC) Tianyuan Fund for Mathematics (Grant No. 12326609), the Guangzhou Key Research and Development Program (Grant No. 2025B03J0125), the Guangzhou Science and Technology Program - Key Research and Development (Grant No.: 2025B03J0155), and the National Natural Science Foundation of China (NSFC) General Program (Grant No.: 82572323).

References

- Walid Al-Dhabyani, Mohammed Gomaa, Hussien Khaled, and Aly Fahmy. 2020. [Dataset of breast ultrasound images](#). *Data in brief*, 28:104863.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millicah, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda

- Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, and 8 others. 2022. Flamingo: a visual language model for few-shot learning. In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22*, Red Hook, NY, USA. Curran Associates Inc.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 others. 2025. *Qwen2.5-vl technical report*. *Preprint*, arXiv:2502.13923.
- Haoran Chen, Junyan Lin, Xinhao Chen, Yue Fan, Xin Jin, Hui Su, Jianfeng Dong, Jinlan Fu, and Xiaoyu Shen. 2025. *Rethinking visual layer selection in multimodal llms*. *Preprint*, arXiv:2504.21447.
- Junying Chen, Chi Gui, Ruyi Ouyang, Anningzhe Gao, Shunian Chen, Guiming Hardy Chen, Xidong Wang, Ruifei Zhang, Zhenyang Cai, Ke Ji, Guangjun Yu, Xiang Wan, and Benyou Wang. 2024. *Huatuogpt-vision, towards injecting medical visual knowledge into multimodal llms at scale*. *Preprint*, arXiv:2406.19280.
- Matthew Christensen, Milos Vukadinovic, Neal Yuan, and David Ouyang. 2024. *Vision-language foundation model for echocardiogram interpretation*. *Nature Medicine*, 30(5):1481–1488.
- Xiaoran Fan, Tao Ji, Changhao Jiang, Shuo Li, Senjie Jin, Sirui Song, Junke Wang, Boyang Hong, Lu Chen, Guodong Zheng, Ming Zhang, Caishuang Huang, Rui Zheng, Zhiheng Xi, Yuhao Zhou, Shihan Dou, Junjie Ye, Hang Yan, Tao Gui, and 5 others. 2024. *Mousi: Poly-visual-expert vision-language models*. *Preprint*, arXiv:2401.17221.
- Haifan Gong, Guanqi Chen, Ranran Wang, Xiang Xie, Mingzhi Mao, Yizhou Yu, Fei Chen, and Guanbin Li. 2021. Multi-task learning for thyroid nodule segmentation with thyroid region prior. In *2021 IEEE 18th international symposium on biomedical imaging (ISBI)*, pages 257–261. IEEE.
- Qiushan Guo, Shalini De Mello, Hongxu Yin, Wonmin Byeon, Ka Chun Cheung, Yizhou Yu, Ping Luo, and Sifei Liu. 2024. *Regiongpt: Towards region understanding vision language model*. *Preprint*, arXiv:2403.02330.
- Xiaoqing Guo, Mohammad Alsharid, He Zhao, Yipei Wang, Jayne Lander, Aris T Papageorghiou, and J Alison Noble. 2026. *A visually grounded language model for fetal ultrasound understanding*. *Nature Biomedical Engineering*, pages 1–17.
- Wenyi Hong, Weihang Wang, Ming Ding, Wenmeng Yu, Qingsong Lv, Yan Wang, Yean Cheng, Shiyu Huang, Junhui Ji, Zhao Xue, Lei Zhao, Zhuoyi Yang, Xiaotao Gu, Xiaohan Zhang, Guanyu Feng, Da Yin, Zihan Wang, Ji Qi, Xixuan Song, and 6 others. 2024. *Cogvlm2: Visual language models for image and video understanding*. *Preprint*, arXiv:2408.16500.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. *Lora: Low-rank adaptation of large language models*. *Preprint*, arXiv:2106.09685.
- Oğuzhan Fatih Kar, Alessio Tonioni, Petra Poklular, Achin Kulshrestha, Amir Zamir, and Federico Tombari. 2024. *Brave: Broadening the visual encoding of vision-language models*. *Preprint*, arXiv:2404.07204.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024a. *Llava-onevision: Easy visual task transfer*. *Preprint*, arXiv:2408.03326.
- Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2023a. *Llava-med: Training a large language-and-vision assistant for biomedicine in one day*. *Preprint*, arXiv:2306.00890.
- Jun Li, Tongkun Su, Baoliang Zhao, Faqin Lv, Qiong Wang, Nassir Navab, Ying Hu, and Zhongliang Jiang. 2024b. *Ultrasound report generation with cross-modality feature alignment via unsupervised guidance*. *IEEE Transactions on Medical Imaging*, 44(1):19–30.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023b. *Blip-2: bootstrapping language-image pre-training with frozen image encoders and large language models*. In *Proceedings of the 40th International Conference on Machine Learning, ICML'23*. JMLR.org.
- Xu Li, Yi Zheng, Haotian Chen, Xiaolei Chen, Yuxuan Liang, Chenghang Lai, Bin Li, and Xiangyang Xue. 2026. *Instruction-guided fusion of multi-layer visual features in large vision-language models*. *Pattern Recognition*, 170:111932.
- Bin Lin, Yang Ye, Bin Zhu, Jiayi Cui, Munan Ning, Peng Jin, and Li Yuan. 2024a. *Video-LLaVA: Learning united visual representation by alignment before projection*. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5971–5984, Miami, Florida, USA. Association for Computational Linguistics.
- Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. 2024b. *VILA: On Pre-training for Visual Language Models*. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 26679–26689, Los Alamitos, CA, USA. IEEE Computer Society.
- Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoub, and Lijuan Wang. 2024a. *Mitigating hallucination in large multi-modal models via robust instruction tuning*. In *The Twelfth International Conference on Learning Representations*.

- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024b. [Improved Baselines with Visual Instruction Tuning](#). In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 26286–26296, Los Alamitos, CA, USA. IEEE Computer Society.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. [Visual instruction tuning](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 34892–34916. Curran Associates, Inc.
- Fadillah Maani, Numan Saeed, Tausifa Saleem, Zaid Farooq, Hussain Alasmawi, Werner Diehl, Ameera Mohammad, Gareth Waring, Saudabi Valappi, Leanne Bricker, and Mohammad Yaqub. 2025. [Fetalclip: A visual-language foundation model for fetal ultrasound image analysis](#). *Preprint*, arXiv:2502.14807.
- Xingqing Nie, Xiaogen Zhou, Tong Tong, Xingtao Lin, Luoyan Wang, Haonan Zheng, Jing Li, Ensheng Xue, Shun Chen, Meijuan Zheng, Cong Chen, and Min Du. 2022. [N-net: A novel dense fully convolutional neural network for thyroid nodule segmentation](#). *Frontiers in Neuroscience*, Volume 16 - 2022.
- David Ouyang, Bryan He, Amirata Ghorbani, Matthew P. Lungren, Euan A. Ashley, David H. Liang, and James Y. Zou. 2019. [Echonet-dynamic: a large new cardiac motion video data resource for medical machine learning](#).
- Lina Pedraza, Carlos Vargas, Fabián Narváez, Oscar Durán, Emma Muñoz, and Eduardo Romero. 2015. [An open access thyroid ultrasound image database](#). In *10th International symposium on medical information processing and analysis*, volume 9287, pages 188–193. SPIE.
- Hanna Piotrkowska-Wróblewska, Katarzyna Dobruch-Sobczak, Michał Byra, and Andrzej Nowicki. 2017. [Open access database of raw ultrasonic signals acquired from malignant and benign breast lesions](#). *Medical Physics*, 44(11):6105–6109.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. [Learning transferable visual models from natural language supervision](#). *Preprint*, arXiv:2103.00020.
- Andrew Sellergren, Sahar Kazemzadeh, Tiam Jaroensri, Atilla Kiraly, Madeleine Traverse, Timo Kohlberger, Shawn Xu, Fayaz Jamil, Cían Hughes, Charles Lau, Justin Chen, Fereshteh Mahvar, Liron Yatziv, Tiffany Chen, Bram Sterling, Stefanie Anna Baby, Susanna Maria Baby, Jeremy Lai, Samuel Schmidgall, and 62 others. 2025. [Medgemma technical report](#). *Preprint*, arXiv:2507.05201.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, and 197 others. 2025a. [Gemma 3 technical report](#). *Preprint*, arXiv:2503.19786.
- LASA Team, Weiwen Xu, Hou Pong Chan, Long Li, Mahani Aljunied, Rui Feng Yuan, Jianyu Wang, Chenghao Xiao, Guizhen Chen, Chaoqun Liu, Zhaodonghui Li, Yu Sun, Junao Shen, Chaojun Wang, Jie Tan, Deli Zhao, Tingyang Xu, Hao Zhang, and Yu Rong. 2025b. [Lingshu: A generalist foundation model for unified multimodal medical understanding and reasoning](#). *Preprint*, arXiv:2506.07044.
- Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. 2024. [Eyes Wide Shut? Exploring the Visual Shortcomings of Multimodal LLMs](#). In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9568–9578, Los Alamitos, CA, USA. IEEE Computer Society.
- Noelia Vallez, Oscar Deniz, Roberto Espinosa, Enrique Santos, and Gloria Bueno. 2025. [Bus-uclm: Breast ultrasound lesion segmentation dataset](#). *Scientific Data*, 12(1):1–10.
- Pavan Kumar Anasosalu Vasu, Hadi Pouransari, Fartash Faghri, Raviteja Vemulapalli, and Oncel Tuzel. 2024. [MobileCLIP: Fast Image-Text Models through Multi-Modal Reinforced Training](#). In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15963–15974, Los Alamitos, CA, USA. IEEE Computer Society.
- Milos Vukadinovic, I-Min Chiu, Xiu Tang, Neal Yuan, Tien-Yu Chen, Paul Cheng, Debiao Li, Susan Cheng, Bryan He, and David Ouyang. 2026. [Comprehensive echocardiogram evaluation with view primed vision language ai](#). *Nature*, 650(8103):970–977.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024a. [Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution](#). *Preprint*, arXiv:2409.12191.
- Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Keqin Chen, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. 2024b. [Cogvlm: Visual expert for pre-trained language models](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 121475–121499. Curran Associates, Inc.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. [Self-instruct: Aligning language models with self-generated instructions](#). *Preprint*, arXiv:2212.10560.

- Beichen Zhang, Pan Zhang, Xiaoyi Dong, Yuhang Zang, and Jiaqi Wang. 2024a. [Long-clip: Unlocking the long-text capability of clip](#). *Preprint*, arXiv:2403.15378.
- Tao Zhang, Xiangtai Li, Hao Fei, Haobo Yuan, Shengqiong Wu, Shunping Ji, Chen Change Loy, and Shuicheng YAN. 2024b. [OMG-LLaVA: Bridging image-level, object-level, pixel-level reasoning and understanding](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Xiaofeng Zhang, Yihao Quan, Chen Shen, Xiaosong Yuan, Shaotian Yan, Liang Xie, Wenxiao Wang, Chaochen Gu, Hao Tang, and Jieping Ye. 2024c. [From redundancy to relevance: Information flow in lvlms across reasoning tasks](#). *Preprint*, arXiv:2406.06579.
- Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong Sun. 2024d. [Llavar: Enhanced visual instruction tuning for text-rich image understanding](#). *Preprint*, arXiv:2306.17107.
- Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. 2025. [Llava-video: Video instruction tuning with synthetic data](#). *Preprint*, arXiv:2410.02713.

A Implementation of Data Generation Pipeline: Desensitization and Quality Control

Data were collected from 15 hospitals across China, encompassing The First Affiliated Hospital of Sun Yat-sen University, The Third Affiliated Hospital of Sun Yat-sen University, The Sixth Affiliated Hospital of Sun Yat-sen University, The Seventh Affiliated Hospital of Sun Yat-sen University, The Eighth Affiliated Hospital of Sun Yat-sen University, The First Affiliated Hospital of Guangxi Medical University, The First Affiliated Hospital of Guangzhou Medical University, Sanshui District People’s Hospital of Foshan City, Puer City People’s Hospital, Sun Yat-sen University Cancer Center, Guangzhou Women and Children’s Health Hospital, Guangzhou First People’s Hospital, The First Affiliated Hospital of Guangzhou University of Chinese Medicine, Guangzhou Military Region General Hospital, and West China Hospital. Given the retrospective nature of this study, all collected data were de-identified prior to the initiation of model training, thereby waiving the requirement for informed consent. Both the data collection protocol and its subsequent utilization were approved by the appropriate institutional ethics review board (IRB).

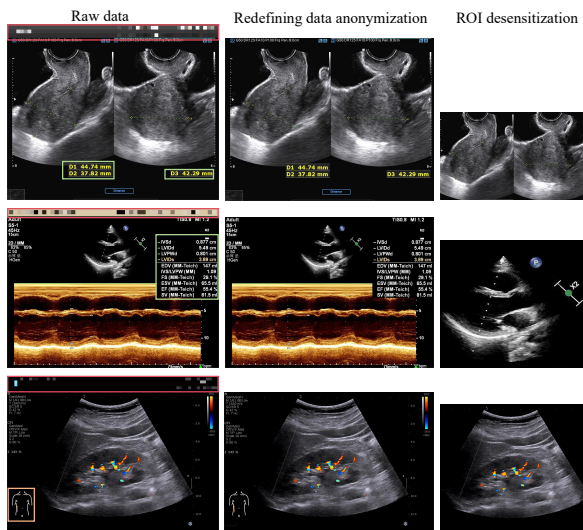


Figure 4: Comparison of traditional and redefined ROI desensitization.

As illustrated in Figure 4, a comparative analysis was conducted between traditional ROI-based image data desensitization and the redefined image data desensitization approach proposed in this study. A core design objective of our redefined

approach is to deliberately retain two types of critical information in the data: meaningful measurement values and anatomical information indicated by anatomical landmarks. Such retention is highly beneficial for improving the performance and clinical applicability of subsequent model training. Specifically, in the leftmost column of raw data, mosaic regions marked with red boxes correspond to sensitive information that needs to be desensitized; green box regions represent the measurement values we aim to preserve, which provide quantitative features for the model and assist it in capturing key numerical information critical for generating accurate and detailed clinical reports; yellow box regions denote the anatomical landmark information we intend to retain, which helps the model accurately locate and understand anatomical regions, thereby enhancing the clinical relevance and interpretability of the model’s outputs.

The data generation pipeline adopts a few-shot prompting strategy to produce high-quality training data for the model, with the utilized opened prompts and closed prompts presented in Figures 5. A diverse array of advanced language and vision-language models was employed for data generation, including GPT-3.5-turbo, GPT-4V, Qwen2-VL-72B, Qwen2.5-VL-7B, Qwen2.5-VL-72B, Llama-3.1-8B, and Llama-3.1-70B. Such diversity in model selection is crucial for enhancing the generalization capability of the subsequent model training process, as it mitigates the risk of overfitting to biased or limited data distributions. Furthermore, strict quality control is integrated into the pipeline to ensure data reliability and clinical validity, with rigorous deduplication during instance creation to preserve content diversity. A newly generated instance is discarded if its ROUGE-L similarity with any previously generated instance exceeds 0.7, or if its Simhash similarity which is quantified by a Hamming distance threshold of ≤ 3 indicates substantial content overlap. For the open-ended data subset, an additional quality control layer is introduced, employing a dual-validation mechanism consisting of sequential AI pre-screening followed by expert verification. The prompts used for AI filtering are illustrated in Figure 6. Following the initial AI pre-screening step, 40,000 samples are randomly selected for clinical validation by domain experts. If errors, inconsistencies, or clinically inappropriate samples are identified in a selected sample set, the generation strategy or the source model associated with that sample set undergoes

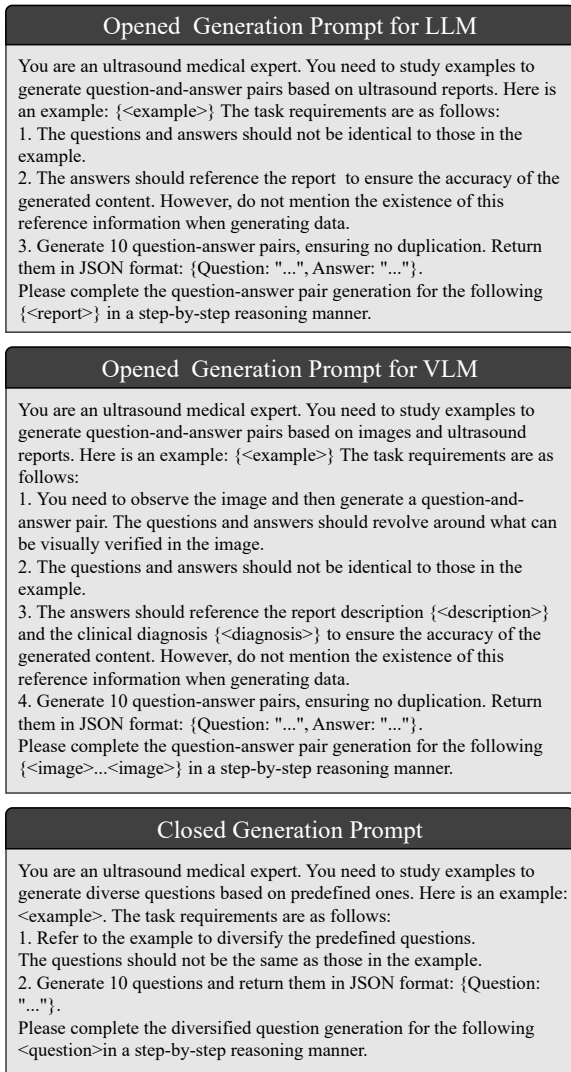


Figure 5: Prompt Templates for Ultrasound Q&A Generation.

retrospective analysis and necessary adjustments, including model replacement and prompt optimization. This entire quality control process ultimately reduces 3.42 million raw instructions to 1.8 million high-quality data pairs.

B Model Architecture and Training Configuration

As presented in Table 9, we detail the model architecture and our proposed two-stage training framework, which integrates a Dual-path MoE for incremental training to enhance domain-specific adaptation. Specifically, the routing experts within this MoE are configured with a dimensionality of 1408, adhering to the Qwen model specifications, while employing a Top-2 selection strategy among four experts to balance specialization and

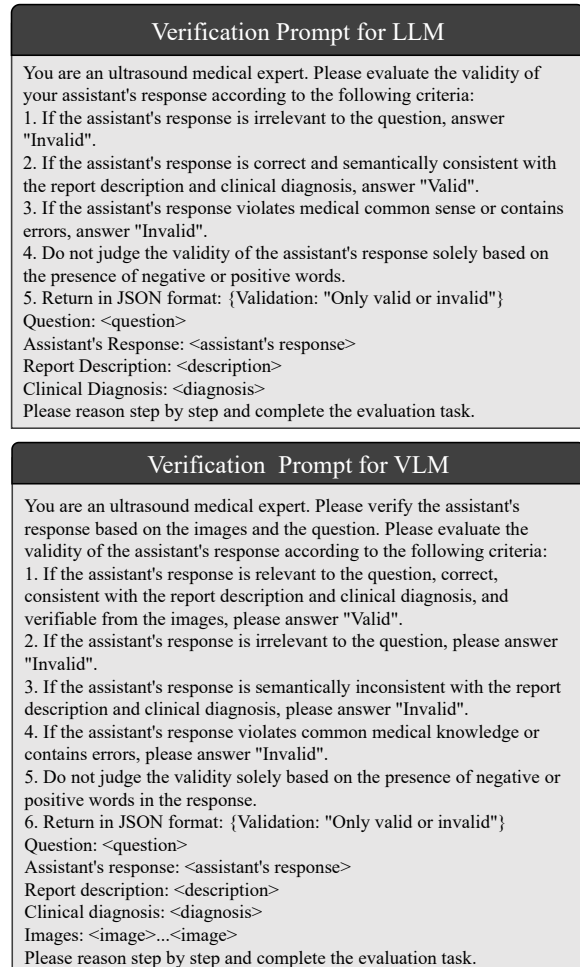


Figure 6: Ultrasound Q&A verification prompt template.

efficiency; concurrently, the shared expert dimension is established at four times that of the routing experts to ensure sufficient capacity for acquiring comprehensive general ultrasound knowledge. Additionally, the PatchMerger module facilitates the fusion and compression of visual representations by consolidating adjacent 2x2 tokens into a single token, thereby optimizing visual feature processing. As shown in Figure 7, during Stage I, the base model parameters are maintained in a frozen state to preserve pre-trained representations, with optimization strictly confined to the newly integrated Dual-path MoE components for acquiring domain-specific visual-textual patterns. In the subsequent Stage II, Cooperative Instruction Fine-tuning is implemented via Low-Rank Adaptation (LoRA), selectively updating the visual encoder, vision-to-text projection layer, and large language model attention layers, while concurrently enabling full-parameter updates for the active ex-

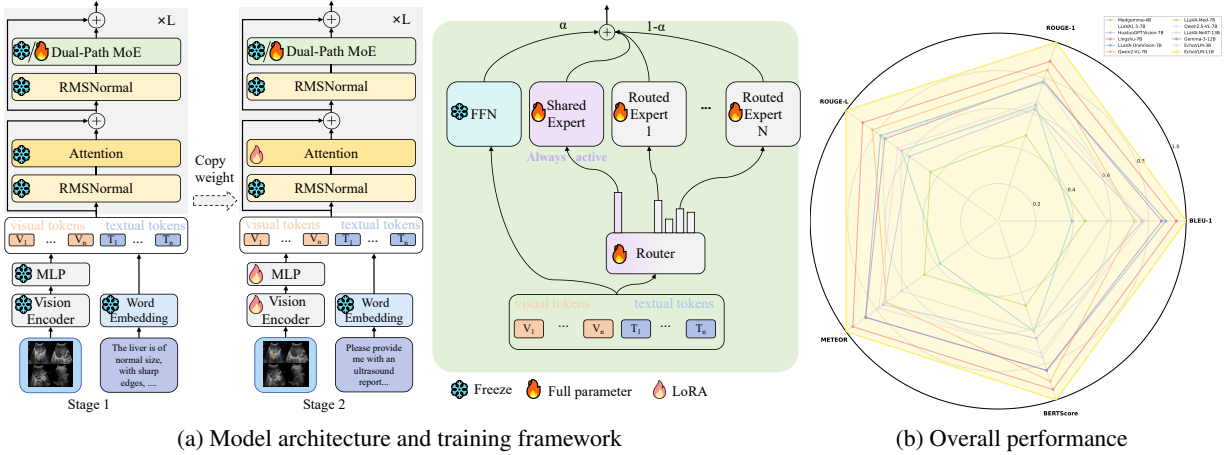


Figure 7: Model structure, training strategy, and overall performance.

Anatomical System	Metric	Medgemma (4B)	LLaVA1.5 (7B)	HuatuoGPT-Vision (7B)	Lingshu (7B)	LLaVA-OneVision (7B)	Qwen2-VL (7B)	LLaVA-Med (7B)	Qwen2.5-VL (7B)	LLaVA-NeXT (13B)	Gemma-3 (12B)	EchoVLM (3B)	EchoVLM (11B)
Breast	Hamming Accuracy	61.22	70.49	72.81	73.03	72.78	74.08	73.38	73.54	72.44	69.59	73.30	75.27
	Macro-Precision	51.09	51.09	62.97	62.11	61.03	62.30	55.77	55.47	61.59	59.70	48.23	59.24
	Macro-Recall	30.81	49.37	59.29	57.27	58.60	57.22	53.05	58.51	52.39	57.66	53.92	57.90
	Macro-F1	35.44	48.07	58.27	56.31	57.16	54.24	49.68	54.55	50.54	57.16	49.18	54.37
Gynecology	Hamming Accuracy	76.42	73.75	75.53	75.71	75.20	75.67	75.67	72.79	75.82	74.24	76.38	78.28
	Macro-Precision	48.98	60.88	45.56	47.52	46.27	42.57	44.05	46.06	65.36	48.56	40.74	56.94
	Macro-Recall	51.13	44.01	52.36	54.25	53.55	46.63	43.99	48.70	44.32	49.85	43.49	56.18
	Macro-F1	47.72	41.67	47.15	48.71	48.22	42.13	39.62	44.36	40.46	47.80	39.60	52.92
Heart	Hamming Accuracy	84.83	77.06	84.28	87.75	85.95	85.31	83.42	86.86	84.56	87.77	85.01	87.84
	Macro-Precision	75.35	86.80	81.95	83.76	80.17	76.70	74.02	78.40	91.48	82.39	73.93	88.60
	Macro-Recall	73.24	67.21	78.35	77.85	82.34	74.69	71.51	80.16	73.32	78.44	74.06	79.31
	Macro-F1	72.23	70.08	78.36	77.81	78.67	74.08	71.21	78.47	74.57	78.24	73.25	80.78
Kidney	Hamming Accuracy	78.93	72.81	75.98	79.01	79.58	74.64	76.28	74.62	76.66	77.92	75.24	81.96
	Macro-Precision	56.03	59.93	58.11	60.73	53.64	53.96	48.42	51.77	60.82	59.86	50.01	65.60
	Macro-Recall	44.53	38.52	48.32	51.51	43.94	38.07	35.38	38.95	36.46	46.78	35.04	61.12
	Macro-F1	45.65	42.01	51.57	53.01	42.24	41.92	36.48	42.19	37.87	49.12	36.95	61.61
Liver	Hamming Accuracy	81.41	70.44	76.73	79.04	79.87	76.67	76.67	71.56	77.04	76.81	82.14	81.19
	Macro-Precision	51.96	53.90	50.24	50.66	51.98	48.56	52.91	48.53	58.34	51.20	48.18	57.22
	Macro-Recall	42.68	38.61	38.61	44.04	48.24	36.53	35.37	31.78	35.88	43.91	43.48	49.77
	Macro-F1	42.32	34.55	41.77	45.96	45.65	39.62	37.83	36.74	38.16	46.13	42.65	50.24
Vessel	Hamming Accuracy	71.89	54.56	61.39	72.22	67.72	68.44	22.28	66.17	54.10	73.06	72.39	51.83
	Macro-Precision	72.85	71.28	79.08	80.22	78.65	79.18	13.61	79.04	69.71	81.70	87.20	92.88
	Macro-Recall	60.20	39.95	55.00	60.34	58.53	50.50	10.05	53.28	39.53	64.10	56.98	33.85
	Macro-F1	63.17	47.94	60.31	63.86	62.86	55.59	4.55	60.35	46.99	68.12	63.22	45.46
Thyroid	Hamming Accuracy	68.81	65.85	75.02	78.50	74.21	78.88	77.45	76.27	76.98	73.85	76.80	80.70
	Macro-Precision	39.82	38.87	47.13	47.49	41.25	54.63	37.06	46.82	36.95	49.42	44.29	54.44
	Macro-Recall	22.74	20.25	36.23	42.31	34.87	39.05	41.41	38.75	41.32	34.50	40.19	47.35
	Macro-F1	27.05	24.35	37.94	41.65	35.73	38.13	36.73	37.26	36.57	38.26	37.43	46.12
Average	Hamming Accuracy	74.79	69.28	74.53	77.89	76.47	76.24	69.31	74.54	73.94	76.18	77.32	76.72
	Macro-Precision	56.58	60.39	60.72	61.78	58.99	59.70	46.55	58.01	63.46	61.83	56.08	67.85
	Macro-Recall	46.48	41.17	52.59	55.37	54.30	48.96	41.54	50.02	46.18	53.61	49.59	55.07
	Macro-F1	47.65	44.10	53.62	55.33	52.93	49.39	39.44	50.56	46.45	54.98	48.90	55.93

Table 8: Multi-label Entity Recognition Performance on Ultrasound Reports.

perts within the Dual-path MoE to maintain their domain-specialized adaptation capabilities. All baseline models in both comparative and preliminary experiments follow the full-parameter fine-tuning paradigm of LLaVA, which ensures comparable parameter updates with EchoVLM across both training stages and enables a fair and rigorous experimental comparison.

C Routing Distributions

As shown in Figure 8, we conduct a comprehensive analysis of the routing distribution within the MoE framework on the test set to rigorously evaluate its efficacy and load-balancing performance. As illustrated in the accompanying figure, the aggregate routing frequencies across all experts exhibit a remarkably balanced pattern, demonstrating that no single expert is disproportionately over- or under-

utilized. Delving deeper into modality-specific routing behavior, we observe that the overwhelming majority of dispatched tokens correspond to image data, whereas textual tokens constitute a markedly smaller fraction. This pronounced imbalance originates from an inherent bias in our dataset composition: although we possess 208,941 radiological reports, these are accompanied by 1.47 million key-frame ultrasound images. Consequently, each individual report is associated with multiple images, thereby skewing the token distribution toward visual content and explaining the dominant image-centric routing observed across experts.

D Visualization

To elucidate the internal decision-making process of VLMs, we employed the Grad-CAM heatmap visualization technique (Zhang et al., 2024c). In

Figure 9, warmer regions (depicted in red) indicate pixels with the strongest gradient flow towards the predicted token logit, signifying that the model primarily attends to these semantically significant regions during textual output generation. Cooler regions (in blue) receive negligible weight allocation, implying minimal contribution to the prediction. The visualized attention patterns confirm that the VLM’s language reasoning is grounded in semantically relevant visual features such as anatomical landmarks, quantitative measurements, and hemodynamic bar charts, as highlighted within red bounding boxes, rather than spurious correlations. However, the model also exhibits attention to non-informative regions including image edges and black areas, as highlighted within blue bounding boxes, which lack semantic utility. This observation highlights a persistent challenge in VLMs: the coexistence of semantic consistency (focusing on task-critical visual elements) and spurious associations (distracting attention to noise or irrelevant artifacts). Future work should enhance model robustness and interpretability through strategies such as attention regularization, constrained attention routing, and explicit mitigation of spurious associations.

E Case Study

Figure 10–16 presents an analysis of report generation cases covering different anatomical regions. Due to space limitations, this section randomly selects examples for discussion. This random selection method ensures an objective assessment of the model’s capabilities while maintaining scientific rigor. The results show that our model exhibits significant advantages over both general-purpose and specialized ultrasound models, effectively capturing clinically meaningful information from images to support accurate reasoning and report generation. Despite these advantages, some limitations exist, including false negatives in nodule identification, suggesting that there is still room for improvement.

F Fine-grained Medical Entity-level Evaluation

To provide a clinically meaningful assessment of generated medical reports, we adopt the entity-level evaluation methodology from Li et al. (2024b). This approach focuses on the accuracy of extracted medical entities rather than surface-level text similarity, which aligns better with how clinicians eval-

Config	Stage I	Stage II
Image encoder	CLIP-ViT-L/14	
Vision-to-text projection	MLP	
LLM	Qwen2-2B/Qwen2-7B	
PatchMerger rate	4	
Shared expert number	1	
Shared expert dimension	1536/5632	
Routing expert number	4	
Routing expert dimension	768/1408	
Top-k	2	
Feature select layer	-1	
Image resolution	392×392	
LoRA rank	-	8
LoRA alpha	-	16
LoRA dropout	-	0.05
Deepspeed	Zero2	
Epoch	1	
Optimizer	AdamW	
Weight decay	0.0	
Learning rate	1e-3	2e-5
Learning rate schedule	Cosine	
Warmup ratio	0.03	
Max length	32768	
Batch size per GPU	1	
GPU	8×A100-80G	
Gradient checkpointing	True	
Precision	Bf16	
Auxiliary loss weight γ	0.001	
Training parameters	1.38B/3.39B	1.39B/3.4B
Total parameters	3B/11B	

Table 9: Model Architecture and Training Configuration.

uate report quality. Table 8 presents fine-grained entity-level performance for ultrasound report generation across seven anatomical systems, while Table 10 shows corresponding results for ultrasound diagnosis. In both tasks, EchoVLM consistently outperforms existing multimodal models across most metrics, demonstrating its superior capability in understanding and modeling fine-grained medical entities.

G Out-of-Distribution (OOD) Evaluation

To perform a more comprehensive assessment of the models, we conduct Out-of-Distribution (OOD) evaluation on the public ultrasound dataset introduced by Li et al. (2024b), with the primary objective of verifying the generalization capability of the evaluated models to unseen ultrasound data. Notably, none of the samples from this public dataset were included in the model training pipeline, thereby eliminating potential data leakage and ensuring the rigor, impartiality, and reliability of the experimental evaluations. To systematically evaluate the OOD generalization capability of all models, we adopt two complementary sets of evalu-

Anatomical	Metric↑	MedGemma (7B)	LLaVA1.5 (7B)	HuatuogPT-Vision (7B)	Lingshu (7B)	LLaVA-OneVision (7B)	Qwen2-VL (7B)	LLaVA-Med (7B)	Qwen2.5-VL (7B)	LLaVA-NeXT (13B)	Gemma-3 (12B)	EchoVLM (3B)	EchoVLM(11B)
Breast	Hamming Accuracy	82.87	79.31	81.75	86.38	82.75	83.62	86.81	86.31	86.34	87.81	88.25	87.94
	Macro-Precision	53.91	48.64	48.12	51.88	44.78	52.85	56.07	65.01	55.25	54.40	55.63	65.43
	Macro-Recall	43.03	42.22	41.48	48.52	45.77	47.56	55.07	51.30	54.45	56.20	57.73	58.11
	Macro-F1	43.63	37.68	41.70	47.39	42.35	46.82	53.29	48.68	52.70	55.11	54.91	57.88
Gynecology	Hamming Accuracy	88.39	87.42	88.53	87.38	86.86	89.28	90.51	87.11	89.59	89.22	90.81	92.52
	Macro-Precision	46.40	62.70	53.17	50.64	48.40	53.52	54.39	44.11	54.63	58.73	59.48	63.53
	Macro-Recall	29.79	22.58	44.03	42.68	41.08	46.02	39.27	43.63	39.42	38.51	50.90	64.28
	Macro-F1	29.86	30.32	41.87	38.72	38.50	43.10	43.12	38.65	43.34	45.23	48.90	62.75
Heart	Hamming Accuracy	87.96	78.29	84.90	90.83	87.67	90.76	93.57	87.90	93.32	90.45	91.58	92.71
	Macro-Precision	57.67	45.12	57.63	70.57	63.84	66.47	72.40	66.73	75.69	65.59	65.13	72.86
	Macro-Recall	50.36	33.47	59.59	60.83	56.19	61.50	65.91	51.52	67.57	62.07	61.93	69.90
	Macro-F1	52.91	37.92	57.82	63.03	58.27	62.27	66.40	54.95	68.47	62.26	61.95	69.93
Kidney	Hamming Accuracy	92.79	81.75	90.78	93.57	94.60	93.73	93.71	92.25	93.71	94.35	94.64	94.54
	Macro-Precision	65.60	39.21	54.89	64.95	60.35	69.22	67.74	57.07	67.73	67.81	72.84	69.66
	Macro-Recall	57.80	31.69	62.58	58.99	59.72	62.37	65.45	58.48	65.49	62.56	64.18	73.58
	Macro-F1	59.11	32.69	56.57	58.33	59.15	63.72	63.87	57.16	63.89	62.51	66.55	70.64
Liver	Hamming Accuracy	92.28	88.09	92.83	93.48	93.50	93.50	92.53	92.57	92.65	90.92	94.24	94.58
	Macro-Precision	39.63	33.09	47.30	47.35	49.63	47.11	45.85	47.21	51.18	39.58	56.50	64.65
	Macro-Recall	43.91	26.64	48.82	46.58	45.82	47.87	47.38	47.61	48.16	44.33	46.04	58.09
	Macro-F1	40.71	26.52	45.16	44.28	46.04	44.74	43.00	43.22	44.33	38.50	47.00	56.35
Vessel	Hamming Accuracy	89.05	84.65	83.40	88.95	87.20	90.30	90.85	88.65	86.14	87.90	91.05	88.20
	Macro-Precision	44.13	32.85	27.12	45.50	45.33	63.67	44.14	37.34	34.75	42.76	43.72	51.06
	Macro-Recall	31.71	28.07	34.30	31.46	34.57	37.03	31.38	29.72	30.67	29.96	32.07	27.19
	Macro-F1	35.40	28.46	28.99	33.83	36.75	44.49	33.55	31.34	30.72	32.65	36.11	34.26
Thyroid	Hamming Accuracy	76.96	69.36	77.31	76.56	79.02	79.14	75.18	78.29	74.87	77.00	82.69	84.78
	Macro-Precision	63.17	54.79	65.49	68.54	66.82	73.65	62.01	68.92	61.57	62.81	71.63	76.73
	Macro-Recall	57.75	33.91	61.29	59.32	64.83	61.22	60.94	62.28	60.38	62.97	68.69	75.05
	Macro-F1	59.25	39.21	59.34	58.93	64.70	60.11	59.69	62.62	59.12	60.54	68.72	71.21
Average	Hamming Accuracy	87.19	81.27	85.64	88.16	87.37	88.62	89.02	87.58	88.30	88.24	90.47	90.75
	Macro-Precision	50.53	45.20	50.53	57.06	54.16	60.93	52.51	55.20	57.26	55.95	60.70	66.27
	Macro-Recall	44.91	31.23	50.30	49.77	49.71	51.94	50.20	49.22	52.31	50.94	54.51	60.89
	Macro-F1	45.84	33.26	47.35	49.22	49.39	52.46	51.85	48.09	51.79	50.69	54.88	60.29

Table 10: Multi-label Entity Recognition Performance on Ultrasound Diagnosis Tasks.

Anatomical	Metric↑	MedGemma (4B)	LLaVA1.5 (7B)	HuatuogPT-Vision (7B)	Lingshu (7B)	LLaVA-OneVision (7B)	Qwen2-VL (7B)	LLaVA-Med (7B)	Qwen2.5-VL (7B)	LLaVA-NeXT (13B)	Gemma-3 (12B)	EchoVLM (3B)	EchoVLM(11B)
Breast	BLEU-1	23.34	19.55	28.45	28.75	26.69	22.55	19.22	23.21	19.82	23.15	22.30	29.16
	ROUGE-1	30.76	26.09	33.88	33.69	32.76	30.36	25.91	31.13	29.35	30.97	30.14	33.96
	ROUGE-L	29.83	23.16	32.18	32.19	31.19	29.37	23.39	30.09	28.08	29.68	29.09	32.35
	METEOR	21.28	18.67	24.55	24.76	23.49	20.75	18.43	21.45	19.32	21.32	20.59	25.01
	BERTScore	52.66	42.09	54.92	54.36	53.96	52.42	42.45	53.16	51.32	52.79	52.24	54.58
Liver	BLEU-1	21.05	23.77	31.95	36.79	18.91	23.73	33.33	22.76	36.31	36.31	23.94	37.36
	ROUGE-1	28.72	27.56	46.03	53.20	29.12	32.89	27.76	49.28	32.38	49.13	35.74	54.06
	ROUGE-L	21.62	21.92	30.01	32.55	21.20	24.70	21.81	31.02	23.81	31.55	24.10	32.85
	METEOR	15.68	17.66	24.34	27.33	14.62	18.63	18.03	25.18	18.62	26.42	18.43	27.68
	BERTScore	43.13	41.66	56.23	61.23	43.76	46.23	41.25	58.72	46.05	58.92	48.72	61.77
Thyroid	BLEU-1	21.48	17.96	27.22	29.84	25.39	19.78	18.59	19.16	15.88	19.72	20.13	30.06
	ROUGE-1	31.55	26.51	35.48	37.27	33.52	31.74	27.05	30.62	28.75	31.13	31.77	37.37
	ROUGE-L	26.07	21.49	28.77	31.57	27.01	25.83	21.54	24.80	23.41	25.34	25.87	31.78
	METEOR	19.71	17.37	23.90	25.86	22.48	18.95	17.50	18.13	15.94	18.80	19.39	26.05
	BERTScore	43.87	37.41	47.78	48.53	46.09	43.98	37.77	43.02	41.13	43.40	43.48	48.61
Average	BLEU-1	21.96	20.43	29.21	31.79	23.66	21.99	20.51	25.23	19.49	26.39	22.12	32.19
	ROUGE-1	30.34	26.72	38.46	41.39	31.80	31.66	26.91	37.01	30.16	37.08	32.55	41.80
	ROUGE-L	25.84	22.19	30.32	32.08	26.47	26.63	22.21	28.64	25.10	28.86	26.39	32.33
	METEOR	18.89	17.90	24.26	25.98	20.20	19.44	17.99	21.59	17.96	22.18	19.47	26.25
	BERTScore	46.55	40.39	52.98	54.71	47.94	47.54	40.49	51.63	46.17	51.70	48.15	54.99

Table 11: Out-of-Distribution Test Performance on Public Datasets Using NLP Metrics.

Anatomical	Metric↑	MedGemma (4B)	LLaVA1.5 (7B)	HuatuogPT-Vision (7B)	Lingshu (7B)	LLaVA-OneVision (7B)	Qwen2-VL (7B)	LLaVA-Med (7B)	Qwen2.5-VL (7B)	LLaVA-NeXT (13B)	Gemma-3 (12B)	EchoVLM (3B)	EchoVLM (11B)
Breast	Hamming Accuracy	35.36	52.89	57.48	56.98	56.55	55.09	52.64	55.58	54.78	55.72	55.09	57.08
	Macro-Precision	37.23	36.98	36.42	36.03	36.13	36.37	38.13	36.77	36.77	36.77	36.77	36.08
	Macro-Recall	45.83	35.21	51.02	46.77	49.17	45.27	37.32	44.35	44.53	46.19	45.08	47.08
	Macro-F1	33.09	29.40	38.31	37.66	36.75	31.89	29.14	33.81	31.23	33.56	31.55	37.87
Liver	Hamming Accuracy	58.04	51.18	76.06	89.61	57.04	56.32	50.73	85.28	52.76	85.97	63.63	91.06
	Macro-Precision	64.44	69.26	58.43	68.54	61.00	63.83	65.78	60.85	59.23	64.16	64.31	58.57
	Macro-Recall	26.48	20.89	45.31	58.75	27.54	24.77	20.35	54.00	21.53	55.72	32.97	60.00
	Macro-F1	34.67	28.25	50.64	58.91	34.84	32.50	27.36	56.19	27.37	57.55	40.78	59.25
Thyroid	Hamming Accuracy	52.79	48.76	58.18	62.03	55.52	52.10	49.20	49.36	46.46	51.30	53.24	62.40
	Macro-Precision	61.97	64.22	60.34	58.10	59.68	61.67	64.89	60.39	57.66	60.18	60.31	57.77
	Macro-Recall	32.33	25.67	40.63	49.38	37.09	30.95	26.41	28.99	24.49	30.31	32.58	49.99
	Macro-F1	33.58	29.90	43.18	49.04	38.11	31.04	30.13	26.90	19.97	30.28	33.88	49.08
Average	Hamming Accuracy	55.40	50.94	63.91	69.54	56.37	54.50	50.86	63.41	51.33	64.33	57.32	70.18
	Macro-Precision	54.55	63.49	51.73	54.22	52.27	53.96	62.93	62.67	48.43	56.74	51.03	50.81
	Macro-Recall	34.88	27.26	45.65	51.63	37.93	33.66	28.03	42.45	30.18	44.07	36.88	52.36
	Macro-F1	33.78	29.18	44.04	48.54	36.57	31.81	28.88	38.97	26.19	40.46	35.40	48.73

Table 12: Out-of-Distribution Test Performance on Public Datasets Using Entity Recognition.

ation metrics, and the detailed OOD test results are presented in the subsequent tables. Table 11 quantifies the OOD performance using conventional natural language processing (NLP) metrics, which assess the coherence, relevance, and accuracy of the generated text. In contrast, Table 12 employs fine-grained entity recognition metrics, which prioritize the precision of medical entity extraction—a critical requirement for clinical ultrasound diagnosis. Collectively, these results demonstrate the models’ adaptability to unseen ultrasound data, a property that is critical for practical clinical ultrasound diagnosis.



Figure 8: Expert routing distribution analysis for different medical imaging tasks: (a) Breast, (b) Gynaecology, (c) Heart, (d) Kidney, (e) Liver, (f) Thyroid, (g) Vessel, and (h) Overall

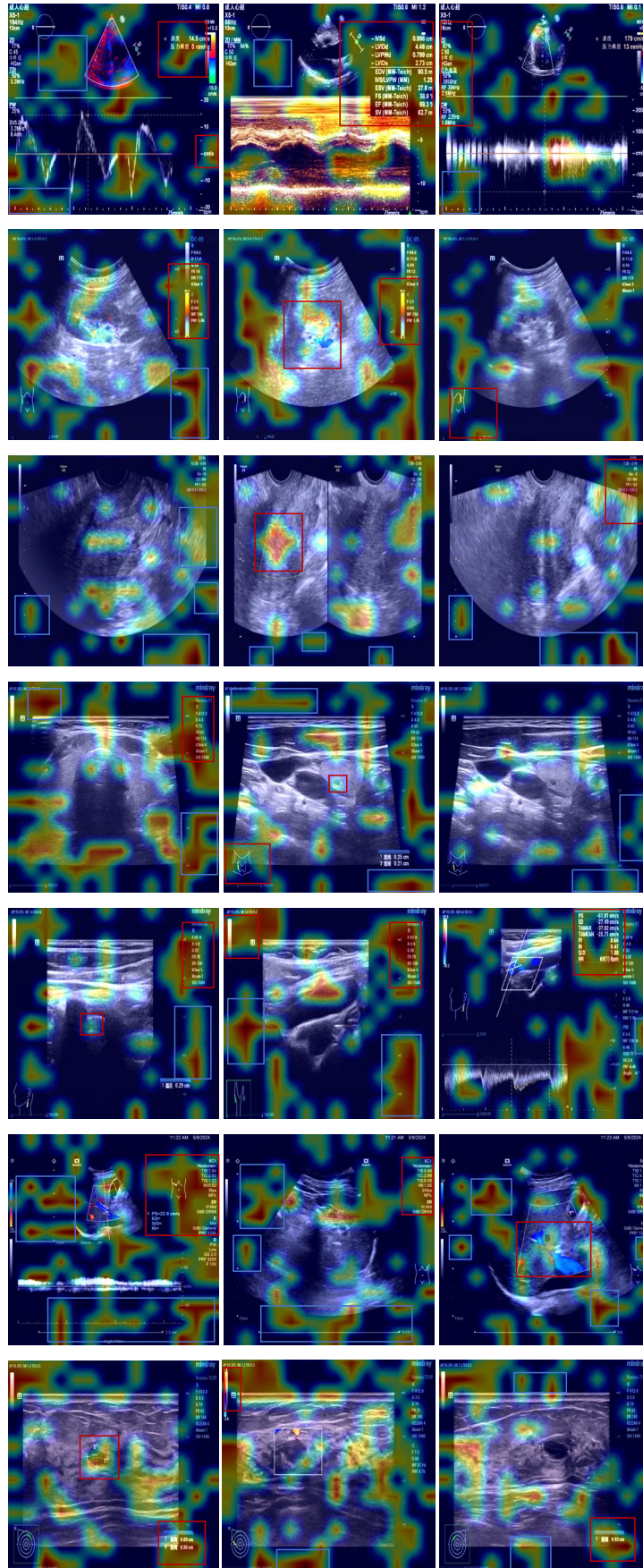


Figure 9: Visualization results. Red boxes highlight meaningful regions, while blue boxes indicate irrelevant or meaningless areas.

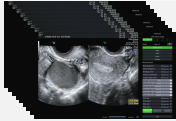
	Can a corresponding medical report be provided based on these ultrasound images?
	
LingShu	Transvaginal, transabdominal, and color Doppler ultrasound examination: The uterus is enlarged, measuring 64mm in length, 52mm in anteroposterior diameter, and 61mm in transverse diameter. The endometrium is 9mm thick and homogeneous. The myometrium shows heterogeneous echogenicity. Hypoechoic areas of 10mm x 8mm are seen on the anterior wall, 11mm x 13mm on the posterior wall, and 14mm x 13mm on the right lateral wall, etc., with a small amount of stellate color Doppler flow signal visible within the lesions. The bilateral ovarian structures are clear, and no obvious abnormal masses are seen in the bilateral adnexa.
Gemma-3	Transvaginal, transabdominal, and color Doppler ultrasound examination: Uterus length 44mm, anteroposterior diameter 31mm, transverse diameter 46mm, normal size. Endometrial thickness 9mm, with slightly heterogeneous echogenicity. Myometrium homogeneous echogenicity. Left ovarian structure clear. No obvious abnormal mass seen in the left adnexal region. Right ovary showed a mass measuring 93mm x 78mm, with clear borders, irregular shape, containing multiple anechoic areas of varying sizes, the largest being 16mm in diameter, with indistinct borders and cloudy echogenicity. A ring-shaped blood flow signal, slightly lighter in color and radiating radially, was observed in conjunction with the posterior cervical canal.
HuatuogPT-Vision	Transvaginal, transabdominal, and color Doppler ultrasound examination: The uterus is enlarged, measuring 64mm in length, 52mm in anteroposterior diameter, and 57mm in transverse diameter. A gestational sac with a maximum diameter of 38mm was observed within the uterine cavity, containing one embryo. The yolk sac was clearly visible, and the embryo's crown-rump diameter was 9mm. Fetal heartbeat was detected on the color Doppler ultrasound. The myometrium showed heterogeneous echogenicity. A hypoechoic area measuring 16mm x 11mm was observed on the anterior wall, containing a small amount of linear color Doppler flow signals. Both ovaries were well-defined, and no obvious abnormal masses were seen in the bilateral adnexa.
LLaVA-OneVision	Transvaginal, transabdominal, and color Doppler ultrasound examination: Uterus length 64mm, anteroposterior diameter 41mm, transverse diameter 48mm, normal size. Endometrial thickness 11mm, homogeneous echo. Myometrium homogeneous echo. Bilateral ovarian structures clear, no obvious abnormal masses seen in bilateral adnexa. Pelvic color Doppler ultrasound showed no obvious abnormal blood flow signals.
Qwen2-VL	Transvaginal, transabdominal, and color Doppler ultrasound examination: Uterus length 54mm, anteroposterior diameter 40mm, transverse diameter 49mm, normal size. Endometrial thickness 10mm, homogeneous echogenicity. Myometrium homogeneous echogenicity. Left ovary structure clear. No obvious abnormal mass seen in the left adnexal region. Right ovary shows an anechoic area of 24mm x 21mm, with clear borders, oval shape, and no color Doppler flow signal seen in the cyst wall.
EchoVLM	Transvaginal, transabdominal, and color Doppler ultrasound examination: Uterus length 64mm, anteroposterior diameter 46mm, transverse diameter 57mm, normal size. Endometrial thickness 10mm, homogeneous echo. Myometrium homogeneous echo. Bilateral ovarian structures clear, no obvious abnormal masses seen in bilateral adnexa. Pelvic color Doppler ultrasound showed no obvious abnormal blood flow signals.
Ground truth	The uterus is enlarged, measuring 64mm in length, 46mm in anteroposterior diameter, and 57mm in transverse diameter. The endometrium is 9mm thick and has heterogeneous echogenicity. The myometrium also shows heterogeneous echogenicity. A 9mm x 8mm hypoechoic area is seen on the left lateral wall. Hypoechoic areas measuring 16mm x 11mm and 14mm x 11mm are seen on the right anterior wall, with a few linear color Doppler flow signals visible within the lesions. The cervical canal structure is present. The left ovary structure is clear. No obvious abnormal mass is seen in the left adnexa. A 44mm x 33mm mass is visible in the right ovary, with clear borders, an oval shape, and cloudy echoes inside; no color Doppler flow signal is seen in the cyst wall.

Figure 10: Gynecological Case Study.

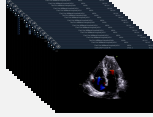
	Can a corresponding medical report be provided based on these ultrasound images?
	
LingShu	Both kidneys are normal in shape and size, with clear outlines and normal parenchymal echogenicity. No focal lesions were observed in the left kidney. Blood flow to both renal hilum is unobstructed. A focal lesion in the right kidney is located in the lower part of the collecting system; it is granular, single, 1.0cm x 0.5cm in size, hyperechoic, with clear borders and acoustic shadowing. The bladder wall is not thickened, the endothelium is smooth, and no lesions were observed within the lumen. No dilation of either ureter was observed. The prostate is normal in size and shape, with homogeneous internal echogenicity, and no lesions were observed.
Gemma-3	Both kidneys are normal in size and shape, with clear outlines. The right renal parenchyma has normal echogenicity, while no space-occupying lesion is seen in the left renal parenchyma. Blood flow is unobstructed at both renal hilum. The bladder wall is not thickened, the endothelium is smooth, and no space-occupying lesion is seen within the lumen. No dilation is seen in either ureter. The prostate measures 4.0cm x 2.4cm x 2.5cm, has a normal shape, homogeneous internal echogenicity, and no space-occupying lesion is seen.
HuatuogPT-Vision	The left kidney measures 12.0cm x 5.3cm with a parenchymal thickness of 0.6cm. The right kidney measures 9.9cm x 4.0cm with a parenchymal thickness of 1.5cm. Both kidneys are normal in shape with clear outlines, and no lesions are seen within the kidneys. Blood flow is unobstructed at both renal hilum. The diameter of the anterior segment of the left renal vein and abdominal aorta is approximately 0.18cm, and the diameter of the distal dilated segment is approximately 0.70cm. The bladder wall is not thickened, the endothelium is smooth, and no lesions are seen within the lumen. Both ureters are not dilated, and no abnormal echoes are seen within the ureters. The prostate is normal in size and shape, with homogeneous internal echoes, and no lesions are seen.
LLaVA-OneVision	The left kidney measures 11.7cm x 5.9cm with a parenchymal thickness of 1.6cm. The right kidney measures 12.3cm x 5.9cm with a parenchymal thickness of 1.0cm. Both kidneys have normal parenchymal echogenicity, clear corticomedullary demarcation, normal collecting system, no focal lesions, and unobstructed blood flow at the renal hilum. The bladder wall is not thickened, the endothelium is smooth, and no lesions are seen within the lumen. Neither ureter is dilated. The prostate is of normal size and shape, with homogeneous internal echogenicity and no lesions are seen.
Qwen2-VL	The liver is of normal size and shape, with a smooth surface and homogeneous parenchymal echogenicity; no space-occupying lesions are observed. Portal vein blood flow is unobstructed and enters the liver. The gallbladder is of normal size and shape, with no lesions found within its lumen. No dilation of intrahepatic or extrahepatic bile ducts is observed. The spleen measures 11.7 cm on its long axis and 5.6 cm in thickness, with no abnormal echogenicity. The pancreas shows no abnormal echogenicity, and the pancreatic duct is not dilated.
EchoVLM	The left kidney measures 11.7cm x 5.6cm with a parenchymal thickness of 1.5cm. The right kidney measures 10.0cm x 4.0cm with a parenchymal thickness of 1.5cm. Both kidneys have normal parenchymal echogenicity, clear corticomedullary demarcation, normal collecting system, no focal lesions, and unobstructed blood flow at the renal hilum. The bladder wall is not thickened, the endothelium is smooth, and no lesions are seen within the lumen. No dilation of either ureter is observed.
Ground truth	The left kidney measures 12.9cm x 5.3cm, with uneven parenchyma thickness, the thickest part being approximately 0.6cm. The collected kidney exhibits irregular separation, resembling a palette, and measures approximately 11.7cm x 5.6cm. The right kidney measures 10.0cm x 4.0cm, with parenchyma thickness of 1.5cm. The right renal parenchyma has normal echogenicity, a clear corticomedullary boundary, a normal collecting system, no focal lesions, and unobstructed blood flow at the renal hilum. No dilation of either ureter is observed.

Figure 11: Kidney Case Study.

<p>I request a clinical report for these ultrasound images.</p>	
<p>LingShu</p>	<p>The liver is of normal size and shape, with a smooth surface and homogeneous parenchymal echogenicity; no space-occupying lesions were observed. Portal vein blood flow was unobstructed, entering the liver with an internal diameter of 0.4 cm and an average flow velocity of 17.9 cm/s. The gallbladder is of normal size and shape, with no lesions found within its lumen. No dilation of intrahepatic or extrahepatic bile ducts was observed. The spleen's long axis is 36.7 cm, and its thickness is 4.1 cm; no abnormal echoes were observed. The pancreas showed no abnormal echoes, and the pancreatic duct was not dilated.</p>
<p>Gemma-3</p>	<p>The liver is of normal size. The liver margins are sharp, the surface is smooth, and the parenchyma shows increased echogenicity. Focal lesion one: located at S5, solitary, measuring 3.7cm × 4.1cm, roughly round in shape, anechoic, with sharp edges, and no clear blood flow signal internally. Focal lesion two: located at S2-S6, multiple, measuring 0.4cm-0.7cm, hyperechoic, roughly round, with visible posterior acoustic shadowing, smooth edges, and no clear blood flow signal. The main portal vein is patent and enters the liver. The average velocity is 17.9cm/s. The gallbladder is of normal shape, and no lesions were found within the lumen. No dilation of intrahepatic or extrahepatic bile ducts was observed. The spleen is of normal size and shows no abnormal echogenicity. The pancreas shows no abnormal echogenicity, and the pancreatic duct is not dilated.</p>
<p>HuatuoGPT-Vision</p>	<p>The liver is of normal size and shape, with a smooth surface and homogeneous parenchymal echogenicity; no space-occupying lesions are observed. Portal vein blood flow is unobstructed, entering the liver with an internal diameter of 0.7 cm and an average flow velocity of 18 cm/s. The gallbladder is of normal size and shape, with no lesions found within its lumen. No dilation of intrahepatic or extrahepatic bile ducts is observed. The spleen's long axis is 36.7 cm, and its thickness is 4.1 cm; no abnormal echoes are seen. The pancreas shows no abnormal echoes, and the pancreatic duct is not dilated.</p>
<p>LLaVA-OneVision</p>	<p>The liver is normal in size and shape, with a smooth surface and homogeneous parenchymal echogenicity; no space-occupying lesions are observed. Portal venous blood flow is unobstructed and enters the liver. The gallbladder is normal in size and shape, with a wall thickness of 0.4 cm. Multiple nodular lesions, ranging in size from 0.4 cm to 0.6 cm, are located in the body of the gallbladder; they are hyperechoic, without posterior acoustic shadowing, and are not mobile. No dilation of intrahepatic or extrahepatic bile ducts is observed. The spleen shows no abnormal echogenicity. The pancreas shows no abnormal echogenicity, and the pancreatic duct is not dilated.</p>
<p>Qwen2-VL</p>	<p>The liver is of normal size and shape, with a smooth surface and homogeneous parenchymal echogenicity; no space-occupying lesions are observed. Portal vein blood flow is unobstructed and enters the liver. The gallbladder is of normal size and shape, with no lesions found within its lumen. No dilation of intrahepatic or extrahepatic bile ducts is observed. The spleen measures 7.6 cm on its long axis and 3.9 cm in thickness, with no abnormal echogenicity. The pancreas shows no abnormal echogenicity, and the pancreatic duct is not dilated.</p>
<p>EchoVLM</p>	<p>The liver is of normal size and shape, with a smooth surface and homogeneous parenchymal echogenicity; no space-occupying lesions were observed. Portal vein blood flow was unobstructed, entering the liver with an internal diameter of 0.7 cm and an average flow velocity of 18.9 cm/s. The gallbladder is of normal size and shape, and no lesions were found within its lumen. No dilation of intrahepatic or extrahepatic bile ducts was observed. The spleen showed no abnormal echogenicity. The pancreas showed no abnormal echogenicity, and the pancreatic duct was not dilated.</p>
<p>Ground truth</p>	<p>The liver is normal in size and shape, with a smooth surface and homogeneous parenchymal echogenicity; no definite space-occupying lesions were observed. The main portal vein has unobstructed blood flow into the liver, with an internal diameter of 0.7 cm and a flow velocity of 18 cm/s. The gallbladder measures 5.5 cm × 2.1 cm, is normal in shape, and has a wall thickness of 0.4 cm. The lesion within the gallbladder is located in the body, is sandy in appearance, and almost completely fills the gallbladder. It is hyperechoic with posterior acoustic shadowing and poor mobility. Heteroechoic areas are seen around the liver, located beside the left lobe and the gallbladder fossa, measuring 4.5 cm × 2.9 cm and 3.7 cm × 4.1 cm respectively, with indistinct borders and sparse blood flow. No dilation of intrahepatic or extrahepatic bile ducts was observed; the common bile duct has an internal diameter of 0.5 cm. The spleen's long axis is 9.0 m, and its thickness is 3.5 cm. No abnormal echogenicity was observed in the spleen. No abnormal echogenicity was observed in the pancreas; the pancreatic duct is not dilated.</p>

Figure 12: Liver Case Study.

Could you provide me with the medical reports for these ultrasound images?



LingShu

Aortic sinus 36mm, ascending aorta 35mm, left atrium 34mm, right ventricle 24mm, interventricular septum 10mm, left ventricle (end-diastole) 46mm, left ventricle (end-systole) 30mm, left ventricular posterior wall 9mm, right atrium 45mm × 30mm, pulmonary artery 24mm. Left ventricular systolic function (Teich method): EF 78%, SV 67ml. Right ventricular systolic function: TAPSE 24mm. Diastolic flow at the mitral valve orifice: E=104cm/s; A=68cm/s; E/A=1.53. Tissue Doppler examination: Mitral valve annulus (interventricular septum): E'=8cm/s; E/E'=13.00; Mitral valve annulus (lateral wall): E'=14cm/s; E/E'=7.43; mean E/E'=10.22; Tricuspid valve anterior annulus: S'=14cm/s. The aortic sinus diameter and ascending aortic diameter are widened. Aortic valve opening is normal; CDFI shows no systolic turbulence and no diastolic regurgitation. Left atrial diameter is normal; no thrombus echoes are seen within the atrial cavity. Mitral valve opening is normal; CDFI detects mild systolic regurgitation. Left ventricular diameter, wall thickness, and wall motion are normal. No definite shunt is seen on atrioventricular level CDFI. Right atrial diameter is normal. Tricuspid valve opening is normal; CDFI detects mild systolic regurgitation. Right ventricular diameter is normal. Aortic and pulmonary artery diameters are normal. Pulmonary valve CDFI detects mild diastolic regurgitation.

Gemma-3

Aortic sinus 31mm, ascending aorta 30mm, left atrium 24mm, right ventricle 20mm, interventricular septum 8mm, left ventricle (end-diastole) 36mm, left ventricle (end-systole) 24mm, left ventricular posterior wall 8mm, right atrium 31mm × 22mm, pulmonary artery 17mm. Left ventricular systolic function (Teich method): EF 74%, SV 49ml. Right ventricular systolic function: TAPSE 17mm. Diastolic flow at the mitral valve orifice: E=54cm/s; A=50cm/s; E/A=1.08. Tissue Doppler examination: mitral valve annulus (interventricular septum): E'=3cm/s; E/E'=18.00; mitral valve annulus (lateral wall): E'=4cm/s; E/E'=13.50; mean E/E'=15.75; tricuspid valve anterior annulus: S'=9cm/s. The diameters of the aortic sinus and ascending aorta are normal. The aortic valve opens normally; CDFI shows no systolic turbulence and no diastolic regurgitation. The left atrial diameter is normal, and no thrombus echoes are seen within the atrial cavity. The mitral valve opens normally; CDFI shows no accelerated diastolic flow and no systolic regurgitation. The left ventricular diameter, wall thickness, and wall motion are normal. No definite shunt is seen on atrioventricular level CDFI. The right atrial diameter is normal. The tricuspid valve opens normally; mild systolic regurgitation is detected on CDFI. The right ventricular diameter is normal. The pulmonary artery diameter is normal. Mild diastolic regurgitation is detected on pulmonary valve CDFI.

HuatuogPT-Vision

Aortic sinus 36mm, ascending aorta 34mm, left atrium 36mm, right ventricle 24mm, interventricular septum 10mm, left ventricle (end-diastole) 45mm, left ventricle (end-systole) 30mm, left ventricular posterior wall 9mm, right atrium 47mm × 30mm, pulmonary artery 24mm. Left ventricular systolic function (Teich method): EF 78%, SV 63ml. Right ventricular systolic function: TAPSE 24mm. Diastolic flow at the mitral valve orifice: E=108cm/s; A=68cm/s; E/A=1.59. Tissue Doppler examination: mitral valve annulus (interventricular septum): E'=8cm/s; E/E'=13.50; mitral valve annulus (lateral wall): E'=14cm/s; E/E'=7.71; mean E/E'=10.61; tricuspid valve anterior annulus: S'=15cm/s. The aortic sinus diameter and ascending aortic diameter are widened. Aortic valve opening is normal; CDFI shows no systolic turbulence and no diastolic regurgitation. Left atrial diameter is enlarged, but no thrombus echoes are seen within the atrial cavity. Mitral valve opening is normal; CDFI detects mild systolic regurgitation. Left ventricular diameter, wall thickness, and wall motion are normal. No definite shunt is seen on atrioventricular level CDFI. Right atrial diameter is normal. Tricuspid valve opening is normal; CDFI detects mild systolic regurgitation, estimated PASP = 30 mmHg. Right ventricular diameter is normal. Aortic and pulmonary artery diameters are normal. Pulmonary valve CDFI detects mild diastolic regurgitation.

LLaVA-OneVision

Aortic sinus 34mm, ascending aorta 29mm, left atrium 35mm, right ventricle 18mm, interventricular septum 10mm, left ventricle (end-diastole) 46mm, left ventricle (end-systole) 24mm, left ventricular posterior wall 7mm, right atrium 44mm × 34mm, pulmonary artery 20mm. Left ventricular systolic function (Teich method): EF 70%, SV 55ml. Right ventricular systolic function: TAPSE 24mm. Diastolic flow at the mitral valve orifice: E=70cm/s; A=60cm/s; E/A=1.17. Tissue Doppler examination: mitral valve annulus (interventricular septum): E'=8cm/s; E/E'=8.75; mitral valve annulus (lateral wall): E'=14cm/s; E/E'=5.00; mean E/E'=6.88; tricuspid valve anterior annulus: S'=14cm/s. The diameters of the aortic sinus and ascending aorta are normal. The aortic valve opens normally; color Doppler flow imaging (CDFI) shows no systolic turbulence or diastolic regurgitation. The left atrial diameter is enlarged, but no thrombus echoes are seen within the atrial cavity. The mitral valve opens normally; CDFI detects mild systolic regurgitation. The left ventricular diameter, wall thickness, and wall motion are normal. No definite shunt is seen on atrioventricular level CDFI. The right atrial diameter is normal. The tricuspid valve opens normally; CDFI detects mild systolic regurgitation, with an estimated PASP of 30 mmHg. The right ventricular diameter is normal. The pulmonary artery diameter is normal. Pulmonary valve CDFI detects mild diastolic regurgitation.

Qwen2-VL

Aortic sinus 36mm, ascending aorta 36mm, left atrium 36mm, right ventricle 24mm, interventricular septum 11mm, left ventricle (end-diastole) 54mm, left ventricle (end-systole) 33mm, left ventricular posterior wall 10mm, right atrium 44mm × 36mm, pulmonary artery 23mm. Left ventricular systolic function (Teich method): EF 71%, SV 103ml. Diastolic flow at the mitral valve orifice: E=87cm/s; A=30cm/s; E/A=2.90. Tissue Doppler examination: mitral annulus (interventricular septum): E'=8cm/s; E/E'=10.88, mitral annulus (lateral wall): E'=10cm/s; E/E'=8.70, mean E/E'=9.79. The diameter of the aortic sinus and ascending aorta is widened. The aortic valve opens normally; CDFI shows no systolic turbulence and no diastolic regurgitation. The left atrial diameter is enlarged, but no thrombus echoes are seen within the atrial cavity. The mitral valve opens normally; CDFI shows no accelerated diastolic flow and no systolic regurgitation. The left ventricular diameter is enlarged, but the wall thickness and wall motion are normal. No definite shunt is seen on atrioventricular level CDFI. The right atrial diameter is normal. The tricuspid valve opens normally; mild systolic regurgitation is detected on CDFI. The right ventricular diameter is normal. The main pulmonary artery diameter is normal. Mild diastolic regurgitation is detected on pulmonary valve CDFI.

EchoVLM

Aortic sinus 36mm, ascending aorta 30mm, left atrium 36mm, right ventricle 24mm, interventricular septum 10mm, left ventricle (end-diastole) 54mm, left ventricle (end-systole) 33mm, left ventricular posterior wall 10mm, right atrium 46mm × 30mm, pulmonary artery 24mm. Left ventricular systolic function (Teich method): EF 71%, SV 94ml. Right ventricular systolic function: TAPSE 24mm. Diastolic flow at the mitral valve orifice: E=78cm/s; A=100cm/s; E/A=0.78. Tissue Doppler examination: mitral valve annulus (interventricular septum): E'=8cm/s; E/E'=9.75; mitral valve annulus (lateral wall): E'=14cm/s; E/E'=5.57; mean E/E'=7.66; tricuspid valve anterior annulus: S'=14cm/s. The aortic sinus diameter is widened, while the ascending aorta diameter is normal. The aortic valve opens normally; color Doppler flow imaging (CDFI) shows no systolic turbulence and no diastolic regurgitation. The left atrial diameter is enlarged, but no thrombus echoes are seen within the atrial cavity. The mitral valve opens normally; CDFI shows no accelerated diastolic flow and no systolic regurgitation. The left ventricular diameter is enlarged, but the wall thickness and wall motion are normal. No definite shunt is seen on atrioventricular level CDFI. The right atrial diameter is normal. The tricuspid valve opens normally; mild systolic regurgitation is detected on CDFI. The right ventricular diameter is normal. The pulmonary artery diameter is normal. Mild diastolic regurgitation is detected on pulmonary valve CDFI.

Ground truth

Aortic sinus 35mm, ascending aorta 36mm, left atrium 42mm, right ventricle 30mm, interventricular septum 11mm, left ventricle (end-diastole) 54mm, left ventricle (end-systole) 32mm, left ventricular posterior wall 10mm, right atrium 46mm × 31mm, pulmonary artery 24mm. Left ventricular systolic function (Teich method): EF 71%, SV 101ml. Right ventricular systolic function: TAPSE 20mm. Diastolic flow at the mitral valve orifice: E=87cm/s; A=30cm/s; E/A=2.90. Tissue Doppler examination: mitral valve annulus (interventricular septum): E'=8cm/s; E/E'=10.88, mitral valve annulus (lateral wall): E'=8cm/s; E/E'=10.88, mean E/E'=10.88, tricuspid valve anterior annulus: S'=14cm/s. The aortic sinus diameter is normal, while the ascending aorta diameter is widened. The aortic valve opens normally; CDFI shows no systolic turbulence and no diastolic regurgitation. The left atrial diameter is enlarged, but no thrombus echoes are seen within the atrial cavity. The mitral valve opens normally; CDFI detects mild systolic regurgitation. The left ventricular diameter is enlarged, but the ventricular wall thickness and motion are normal. A left-to-right shunt is detected at the foramen ovale in the atrial septum via CDFI, with a shunt width of approximately 3.2 mm. No definite shunt is detected at the ventricular level via CDFI. The right atrial diameter is normal. The tricuspid valve opens normally; CDFI detects mild systolic regurgitation, with an estimated PASP of 34 mmHg. The right ventricular diameter and anterior wall thickness are normal. The pulmonary artery diameter is normal. Mild diastolic regurgitation is detected at the pulmonary valve via CDFI. The inferior vena cava diameter is normal, with an inspiratory collapse rate greater than 50%.

Figure 13: Heart Case Study.


	 <p>Could you provide me with the analysis reports for these ultrasound images?</p>
LingShu	The intima of both carotid arteries was not smooth, with the thickest intima-media thickness (IMT) being 1.2 mm. Blood flow in both internal jugular veins and vertebral veins was patent, and no thrombosis was observed.
Gemma-3	The intima of both carotid arteries was smooth, with the thickest intima-media thickness (IMT) being 0.8 mm. Internal carotid artery PSV: left 47 cm/s, RI 0.71; right 44 cm/s, RI 0.69. Blood flow in both internal jugular veins and vertebral veins was patent, with no thrombosis observed.
HuatuoGPT-Vision	The intima of both carotid arteries was not smooth, with the thickest intima-media thickness (IMT) being 1.2 mm. Internal carotid artery PSV: Left 44 cm/s, RI 0.69; Right 47 cm/s, RI 0.71. Blood flow in both internal jugular veins and vertebral veins was patent, with no thrombosis observed.
LLaVA-OneVision	The intima of both carotid arteries was not smooth, with the thickest intima-media thickness (IMT) being 1.4 mm. Blood flow in both internal jugular veins and vertebral veins was patent, and no thrombosis was observed.
Qwen2-VL	The intima of both carotid arteries was smooth, with the thickest intima-media thickness (IMT) being 0.8 mm. Blood flow in both internal jugular veins and vertebral veins was unobstructed, and no thrombosis was observed.
EchoVLM	The intima of both carotid arteries was smooth, with the thickest intima-media thickness (IMT) being 0.8 mm. Internal carotid artery PSV: left 47 cm/s, RI 0.71; right 44 cm/s, RI 0.69. Blood flow in both internal jugular veins and vertebral veins was patent, with no thrombosis observed.
Ground truth	The intima of both carotid arteries (common carotid artery, internal carotid artery, and external carotid artery) was smooth, with a maximum intima-media thickness (IMT) of 1.1 mm and normal luminal diameter. Blood flow in both internal jugular veins and vertebral veins was unobstructed, and no thrombosis was observed.

Figure 14: Vascular Case Study.

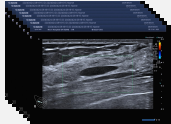
	 <p>Rate these ultrasound images based on the BI-RADS scoring system.</p>
LingShu	birads 2
Gemma-3	birads 3
HuatuoGPT-Vision	Bilateral breast lesions, suggestive of fibroadenoma (BI-RADS category 3). Bilateral cystic hyperplasia of the breast. No abnormally enlarged lymph nodes were found in either axilla.
LLaVA-OneVision	birads 3
Qwen2-VL	birads 2
EchoVLM	birads 3
Ground truth	birads 3

Figure 15: Breast Case Study.


	 <p>Please issue a detailed medical report based on the provided ultrasound images.</p>
LingShu	The left lobe of the thyroid gland is normal in size and shape, with coarse and mixed echogenicity. It presents a single focal lesion measuring 1.6cm × 0.7cm, with mixed echogenicity, liquefaction, no calcification, no halo, indistinct borders, and sparse blood supply. The right lobe of the thyroid gland is normal in size and shape, with coarse and mixed echogenicity. It presents multiple focal lesions, ranging from 0.2cm to 0.8cm in size, with mixed echogenicity. Some nodules show liquefaction, no calcification, no halo, indistinct borders, and sparse blood supply. The isthmus of the thyroid gland is 0.3cm thick. No abnormally enlarged lymph nodes were observed in either side of the neck.
Gemma-3	The bilateral thyroid lobes and isthmus are of normal size and regular shape, with a smooth capsule and coarse internal echoes, but no obvious abnormal echoes were seen, and the blood supply is normal. No abnormally enlarged lymph nodes were seen in either side of the neck.
HuatuoGPT-Vision	The left lobe of the thyroid gland is normal in size and shape, with coarse and mixed echogenicity. Multiple focal lesions, ranging from 0.4cm to 1.6cm in size, are present, with mixed echogenicity. Some nodules show liquefaction, but no calcification, no acoustic halo, indistinct borders, and sparse blood supply. The right lobe of the thyroid gland is normal in size and shape, with coarse and mixed echogenicity. Multiple focal lesions, ranging from 0.5cm to 2.7cm in size, are present, with mixed echogenicity. Some nodules show liquefaction, but no calcification, no acoustic halo, indistinct borders, and sparse blood supply. The isthmus of the thyroid gland is 0.3cm thick. No abnormally enlarged lymph nodes were observed in either side of the neck.
LLaVA-OneVision	The bilateral breast tissue showed coarse and mixed echogenicity, with uniform dilation of the mammary ducts and normal blood flow. The left lesion was a single, hypochoic mass located at the 12 o'clock position, measuring 0.7cm × 0.4cm, with clear borders and a capsule. The right lesion was a single, hypochoic mass located at the 9 o'clock position, measuring 0.8cm × 0.5cm, with clear borders and a capsule. No abnormally enlarged lymph nodes were observed in either axilla.
Qwen2-VL	The left lobe of the thyroid gland is normal in size and shape, with coarse and mixed echogenicity. Multiple focal lesions, ranging from 0.3cm to 1.0cm in size, are present, with mixed echogenicity. Some nodules show liquefaction, but no calcification, no acoustic halo, indistinct borders, and sparse blood supply. The right lobe of the thyroid gland is normal in size and shape, with coarse and mixed echogenicity. Multiple focal lesions, ranging from 0.3cm to 1.0cm in size, are present, with mixed echogenicity. Some nodules show liquefaction, but no calcification, no acoustic halo, indistinct borders, and sparse blood supply. The isthmus of the thyroid gland is 0.3cm thick. No abnormally enlarged lymph nodes were observed in either side of the neck.
EchoVLM	The left lobe of the thyroid gland is normal in size and shape, with coarse and mixed echogenicity. It presents a single focal lesion, measuring 0.7cm × 0.6cm, with mixed echogenicity, liquefaction, no calcification, no halo, indistinct borders, and sparse blood supply. The right lobe of the thyroid gland is normal in size and shape, with coarse and mixed echogenicity. It presents multiple focal lesions, ranging from 0.3cm to 1.0cm in size, with mixed echogenicity. Some nodules show liquefaction, no calcification, no halo, indistinct borders, and sparse blood supply. The isthmus of the thyroid gland is 0.2cm thick. No abnormally enlarged lymph nodes were observed in either side of the neck.
Ground truth	The right lobe of the thyroid gland is normal in size and shape, with coarse and mixed echogenicity. A single focal lesion, measuring 1cm × 0.56cm × 0.70cm, with an aspect ratio >1, is irregular in shape, hypochoic, without liquefaction, and contains scattered punctate hyperechoic areas. There is no halo, the borders are clear, and the blood supply is rich. The left lobe of the thyroid gland is normal in size and shape, with coarse and mixed echogenicity, and no focal lesions are observed. The isthmus of the thyroid gland is 0.3cm thick. No abnormally enlarged lymph nodes are seen in either side of the neck.

Figure 16: Thyroid Case Study.