

EPISTEMIC-AWARE VISION-LANGUAGE FOUNDATION MODEL FOR FETAL ULTRASOUND INTERPRETATION

Anonymous authors
 Paper under double-blind review

ABSTRACT

Recent medical vision-language models have shown promise on tasks such as VQA, report generation, and anomaly detection. However, most are adapted to structured adult imaging and underperform in fetal ultrasound, which poses challenges of multi-view image reasoning, numerous diseases, and image diversity. To bridge this gap, we introduce **FetalMind**, a medical AI system tailored to fetal ultrasound for both report generation and diagnosis. Guided by clinical workflow, we propose *Salient Epistemic Disentanglement* (SED), which injects an expert-curated bipartite graph into the model to decouple view-disease associations and to steer preference selection along clinically faithful steps via reinforcement learning. This design mitigates variability across diseases and heterogeneity across views, reducing learning bottlenecks while aligning the model’s inference with obstetric practice. To train **FetalMind** at scale, we curate **FetalSigma-1M** dataset, the first large-scale fetal ultrasound report corpus, comprising 20K reports from twelve medical centers, addressing the scarcity of domain data. Extensive experiments show that **FetalMind** outperforms open- and closed-source baselines across all gestational stages, achieving +14% average gains and +61.2% higher accuracy on critical conditions while remaining efficient, stable, and scalable.

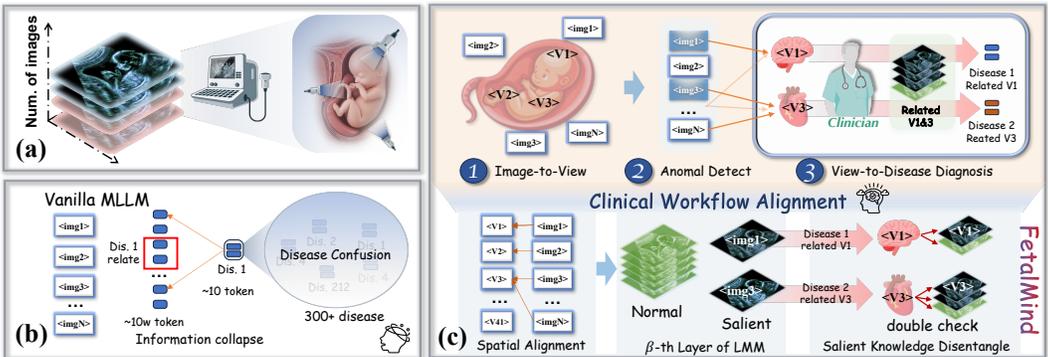


Figure 1: (a) Fetal ultrasound workflow. (b) Limitations of vanilla MLLMs on multi-view scans: ❶ A severe imbalance, with abundant visual tokens but limited textual supervision, induces representation collapse; ❷ Fetal imaging spans > 300 fine-grained diseases, markedly complicating robust diagnosis. (c) **FetalMind** aligns with the clinical workflow: view examination, abnormality detection, and disease tracing via knowledge.

1 INTRODUCTION

Ultrasound is the preferred tool for prenatal assessment, routinely used to track fetal growth, monitor pregnancy progression, and support clinical diagnosis (Salomon et al., 2022; Neilson et al., 1996). In contrast to adult imaging, fetal ultrasound requires integrating information across multiple views and gestational stages (Azad et al., 2024). Effective diagnosis must jointly consider developmental trajectories and early indicators of potential abnormalities (Lee et al., 2023). As illustrated in Figure 1a, fetal ultrasound typically involves many images with inconsistent view counts, substantial inter-case heterogeneity, and pronounced disease variability (Krishna & Kokil, 2024).

With the rise of deep learning, prior satisfactory works has decomposed fetal ultrasound into sub-tasks, e.g., biometric measurement, view classification, gestational age estimation, and anomaly

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

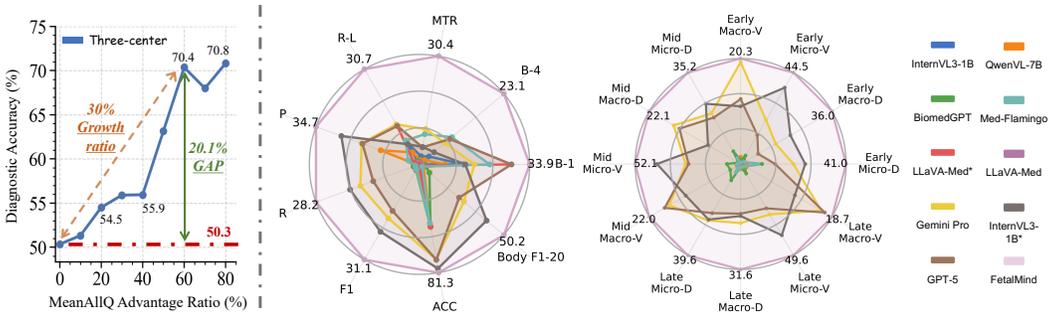


Figure 2: *Left:* Positive correlation (> 0.3) between diagnostic accuracy and the relative attention advantage of disease-related over non-disease views. Attention is measured by MeanALLQ, defined as the mean attention weight over all query tokens across layers and heads, and results are shown for Qwen-VL 2.5. *Right:* Multi-center evaluation of report generation and diagnosis with trimester-level diagnostic performance comparison.

analysis, achieving encouraging task-specific results (Fiorentino et al., 2023). More recently, several outstanding medical MLLM models have been proposed to handle cross-modal medical image and text instruction tasks, demonstrating significant results in experiments (Moor et al., 2023a).

However, when aligning multiple images with text, existing medical MLLM exhibit two critical issue (see Figure 1b): ❶ **Information collapse.** During disease-image alignment, diagnosis often contain only ~ 10 text tokens, while the associated image evidence may expand to $\sim 10^4$ visual tokens across views; the severe imbalance causes salient cues to be drowned out or ignored. ❷ **Disease confusion.** Fetuses present with multiple coexisting conditions, and disease-relevant views frequently overlap or partially align across slices. Such complexity hinders the inter-disease discriminability and results in confounded anomaly recognition and diagnosis. Consequently, reliable fetal ultrasound report generation and diagnosis remain unachieved with current deep learning approaches, limiting both clinical automation and decision support (Slimani et al., 2023).

The core challenge arises from the limitations of current MLLM approaches, which remain constrained to single-image, image-text alignment and therefore fail to capture anatomical development and latent abnormality associations across multiple views (Cheng et al., 2025; Liu et al.). In clinical practice, however, fetal ultrasound diagnosis does not rely on isolated images; it integrates spatial continuity and the developmental logic of anatomy across views (Carvalho et al., 2023). Existing models, lacking the ability to disentangle complementary information across views, often blur the correspondence between views and disease features (Arnaout et al., 2021). As illustrated in Figure 2 left, insufficient attention to disease-relevant views frequently leads to hallucinated or biased diagnoses, undermining reliability and diverging from established clinical workflows. In contrast, as illustrated in Figure 1c, obstetricians begin with a comprehensive survey of all views and progressively refine their focus on multiple views of specific regions to ensure thorough assessment.

Motivated by clinical workflows, we introduce *Spatial Alignment* to capture image-to-view correspondences and integrate it with *Salient Epistemic Disentanglement* through view preference optimization (SVPO). This synergy enhances the model’s sensitivity to disease-bearing planes while explicitly injecting disease-plane associations, enabling the joint disentanglement of salient versus normal planes at both the case and view levels. Such modeling mirrors the reasoning process of obstetricians (Figure 1c), steering inference toward clinically grounded, auditable, and verifiable reports, thereby avoiding “isolated image \rightarrow conclusion” shortcuts. To train **FetalMind** effectively, we construct the first large-scale fetal ultrasound report dataset, **FetalSigma-1M**. The dataset consists of real-world clinical data collected from 12 medical centers, covering 20,566 patients with 1.19M ultrasound images paired with expert-verified reports and diagnoses across early, mid, and late trimesters. As shown in Figure 2 right, **FetalMind** surpasses state-of-the-art medical MLLMs and general-purpose MLLMs (e.g., GPT-5) across multiple downstream tasks, highlighting its robustness and clinical applicability. To summarize, our contributions as follows:

- ❶ To the best of our knowledge, we present **FetalMind**, the first model for fetal ultrasound report generation and diagnosis capable of handling a variable number of views, with **1B** and **7B** versions. The model integrates salient epistemic disentanglement with salient view preference optimization and bipartite knowledge graph to capture disease-view associations, explicitly decouple salient from normal views at both the disease and view levels.

- ② We construct **FetalSigma-1M**, a large-scale multi-center benchmark comprising 1M multi-view ultrasound images and 20K paired clinical reports. The dataset spans all trimesters, covers all standard views, and includes over 300 diseases categories derived from real clinical examinations.
- ③ We conduct extensive experiments showing that **FetalMind** achieves a 14% improvement in multi-center and zero-shot multi-device diagnosis, while maintaining strong robustness and generalization across diverse real-world clinical scenarios.

2 RELATED WORK

Medical Multimodal Large Language Models. Building on the success of general multimodal large language models (MLLMs) such as CLIP (Radford et al., 2021) and GPT-4 (Achiam et al., 2023), recent efforts have explored foundation models for medicine that learn unified image-text representations. LLaVA-Med augments biomedical imagery with open-ended dialogue and QA via large-scale chart-caption data and GPT-4-based instruction synthesis (Li et al., 2023). Med-PaLM accommodates text, images, and genomics under a single parameterization (Singhal et al., 2025). Several medical MLLMs also incorporate ultrasound data. BiomedGPT is an open, lightweight medical VLM supporting images, text, and tables (Zhang et al., 2024). HealthGPT unifies multi-modal understanding and generation in an autoregressive framework (Lin et al., 2025). MedRegA provides a bilingual, general-purpose medical AI across eight modalities for both image- and region-level vision-language tasks (Wang et al., 2024). As a general foundation model, GPT-5, exhibits strong cross-modal reasoning and, with instruction tuning and domain adaptation, can support medical VQA, report generation, and clinical decision support (Hou et al., 2025). Despite this progress, most prior work targets adult images, with limited coverage of obstetrics and fetal ultrasound, which is a basic tool for prenatal care. Multi-center heterogeneity and the complexity of multi-image/multi-view inputs remain open challenges. Overall, existing methods remain task-specific and confined to per-view analysis, whereas clinical practice requires aggregating information across multiple views to support diagnosis and decision-making. To the best of our knowledge, no existing AI model and dataset specifically address fetal ultrasound report generation and diagnosis.

Fetal Ultrasound. Ultrasound is the primary imaging modality for fetal anomaly screening, yet substantial appearance variability, scale differences, disease diversity, and multi-view images make automated interpretation challenging (Hu et al., 2023). Prior work has largely relied on supervised learning on single views, emphasizing standard-view recognition and automated biometry (Awadalla et al., 2023). In multi-image MLLM studies, Liu et al. employ DPO to guide models to attend to text-relevant regions across multiple images; however, these images often lack intrinsic inter-image dependencies. FetalCLIP learns anatomy-sensitive, generalizable representations via large-scale text-image contrastive learning and cross-modal alignment, benefiting downstream tasks such as classification and gestational-age estimation (Maani et al., 2025). *The aforementioned works, e.g., FetalCLIP, operate at the level of single-image parsing within the clinical workflow (Figure 1c, (I) Image-to-View), focusing on view classification, organ segmentation, and a limited set of related sub-tasks, primarily to assist clinicians in identifying standard views. Beyond this, FetalMind is the first to achieve holistic interpretation of fetal ultrasound images and can directly generate full reports and diagnostic conclusions that support clinical decision-making.*

3 CLINICAL FETAL ULTRASOUND DATASET CONSTRUCTION

In this section, we introduce the **FetalSigma-1M** dataset, composed of three subsets: ① Image-Report dataset: 20K image-report pairs, where each case includes multiple ultrasound images and a fine-grained clinical report covering biometric measurements, structural assessments, and abnormal findings. ② Image-Diagnosis dataset: 1M images organized as multi-image, case-level samples paired with physician-verified diagnostic reports. ③ View Classification dataset: 10K fetal ultrasound images with fine-grained view annotations collected across three medical centers.

3.1 IMAGE-REPORT DATASET

Scope & Scale. We curate a large-scale, multi-center dataset for fetal ultrasound report generation and disease diagnosis that spans the full gestational spectrum and all fetal systems. The cohort comprises Early 5.0K, Mid 10.9K, and Late 5.2K examinations. Class balance is maintained with

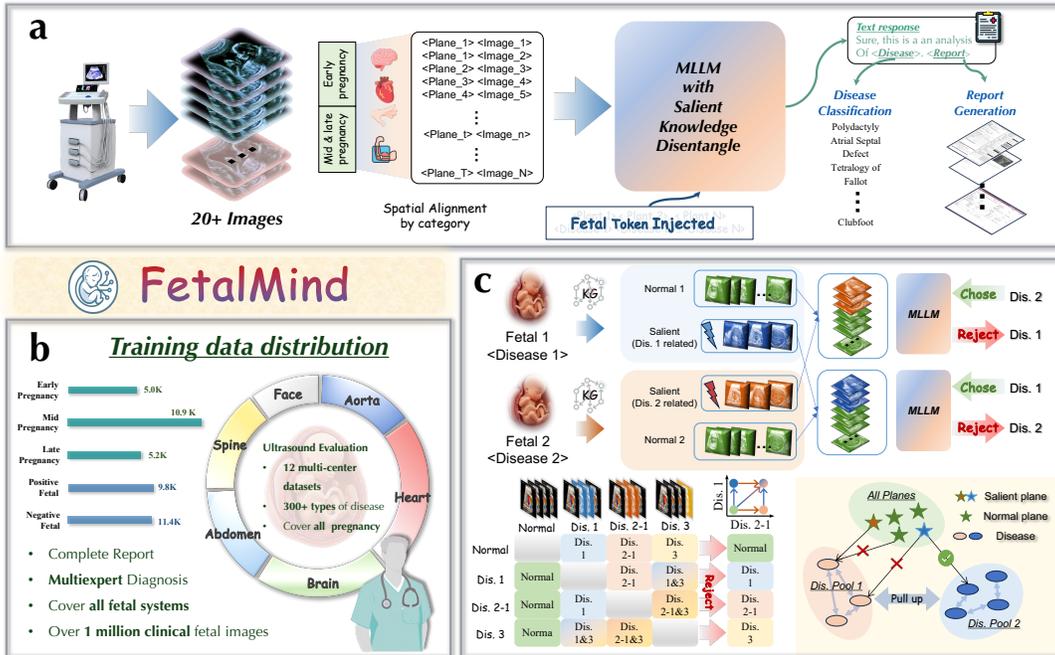


Figure 3: (a) *FetalMind* aligns with clinical cognition by classifying images into pregnancy-specific views, encoding disease-view keywords as special tokens, and reinforcing their intrinsic associations via salient epistemic disentanglement (SED). (b) *FetalSigma-1M* comprises 1 million fetal image-report-diagnosis triplets in 12 centers. (c) Overview of SED. Salient views are identified from disease-view graphs (see Section A.7) and treated as perturbation variables, swapped across fetuses with disease replaced. Bottom-left: intersection- and union-based substitution between diseased regions and views. Bottom-right: SVPO not only injects disease-view knowledge graphs into MLLMs but also enhances inter-disease discriminability

9.8K positive and 11.4K negative cases across 300+ disease categories. Data originate from 12 centers and multiple device models, totaling > 1M clinical ultrasound images and enabling robust evaluation of cross-center generalization. Structured documentation across the heart, central nervous system, chest, abdomen, spine, face, neck, and long bones, covering all fetal systems, to support fine-grained fetal ultrasound analysis and multi-image modeling.

Curation & Splits. We apply unified multi-center cleaning, de-duplication, and quality control, including removal of low-quality frames and harmonization across devices/exports. Our survey (see Figure 1a) indicates that medical MLLMs trained on generic image-text pairs frequently miss diagnoses, which is an unacceptable failure mode in clinical practice. Accordingly, during curation we deliberately enriched positive cases to stabilize supervision, as routine fetal screening exhibits a base positive rate of < 1% in our observations across more than three centers. All positive case reports were finalized under the diagnoses of at least two expert clinicians.

3.2 IMAGE-DIAGNOSIS DATASET

Because many reports lack explicit diagnostic statements, we assigned a *Diagnosis* to each examination under physician supervision. Specifically, we constructed a disease ontology with 310 entities and their corresponding anatomical sites. Each report was then processed with DeepSeek-R1 (Guo et al., 2025) to extract provisional diagnoses by referencing this ontology, after which multiple expert fetal sonographers reviewed and corrected the outputs to obtain finalized diagnoses.

View Classification dataset. Accurate view localization from ultrasound video frames is the first step in fetal examination, as subsequent measurements and diagnoses rely on the correct anatomical view. Guidelines require nearly 20 standard views in the second trimester, with substantial variation across gestational stages and fetal positions, making automated modeling challenging (Salomon et al., 2022). To ensure reliable supervision, we annotated 11,358 images from three centers in *FetalSigma-1M* into 54 view categories under expert guidance, covering early, mid, and late gestation and including key views such as four-chamber, aortic arch, and three-vessel views. This subset is used to train the multi-view classification model in the Spatial Alignment stage (see Figure 1c).

4 METHODOLOGY

Figure 3a outlines how **FetalMind** is deployed within a fetal-ultrasound pipeline. Guided by clinical workflow, given multiple input images, **FetalMind** first performs *spatial alignment* to map each image to its anatomical view (\triangleright Section 4.1), followed by *fetal token injection* to encode domain priors and mitigate disease confusion induced by text similarity (\triangleright Section 4.2). We then describe how *view-disease swapping* constructs positive/negative pairs and how SVPO strengthens the model’s preference for disease-relevant views (\triangleright Section 4.3). Finally, we present the principles of multi-view swapping under different conditions. Please refer to Section B for more details.

4.1 CLASS-WISE SPATIAL ALIGNMENT

Identifying the correct imaging view is a prerequisite for reliable fetal diagnosis and report generation. To align with the view-image paradigm and remain robust against imaging noise, fetal pose variation, and gestational-age differences, we adopt a classification-based strategy. Given the substantial distribution shift between early, mid, and late gestation, and the clinical practice of treating them as distinct tasks, we partition the 10K view-annotated images in **FetalSigma-1M** into *early* and *mid/late* subsets, using a 7:1:2 train/val/test split for pretraining. As illustrated in Figure 3a, the spatial alignment module incorporates two classifiers (Woo et al., 2023), trained separately on the double model. The early-gestation model spans 9 views categories, while the mid- and late-gestation model covers 41 categories, encompassing all clinically essential planes (Pellerito et al., 2018).

4.2 FETAL TOKEN INJECTION

We introduce the *Fetal Token Injection* strategy to explicitly encode domain-specific priors in fetal ultrasound by mapping key terms to spatial tokens. The rationale stems from the holistic nature of the fetus: although over 300 congenital anomalies have been documented, many exhibit highly similar linguistic descriptions (e.g., *ventricular septal defect* vs. *atrial septal defect*), yet correspond to clinically distinct diseases with divergent prognoses and management strategies. Similarly, prenatal ultrasound defines more than 40 standard imaging planes. While their textual descriptions may partially overlap, these planes are not interchangeable in clinical workflows. Without explicit token-level disentanglement, MLLMs tend to conflate semantically similar but clinically independent entities, ultimately yielding unreliable predictions and hallucinated report generation. This strategy introduces structured, view- and disease-aware tokens that enforce clear separability among near-synonymous terms and imaging planes, enhancing the reliability of diagnosis and reporting.

4.3 SALIENT EPISTEMIC DISENTANGLEMENT

Each fetus i is represented as a multi-view sample $\mathcal{X}^i = (p, I_{i,p}), p \in \mathcal{P}$, where \mathcal{P} denotes the set of anatomical views and $I_{i,p}$ the image for view p . View-image correspondence $(p, I_{i,p})$ is obtained by the class-wise spatial alignment (Section 4.1). As shown in Figure 3c, the clinically confirmed disease set is $\mathcal{D}_i \subseteq \mathcal{V}_{\text{dis}}$. We construct an expert-curated disease \rightarrow view *bipartite knowledge graph* $G : \mathcal{V}_{\text{dis}} \rightarrow 2^{\mathcal{P}}$ under the guidance of textbooks and experts that maps each disease d to its salient views $G(d) \subseteq \mathcal{P}$. Given d , define the salient and non-salient view sets $\mathcal{P}^{(+)}(d) = G(d)$, $\mathcal{P}^{(-)}(d) = \mathcal{P} \setminus G(d)$ and split \mathcal{X}_i as $\mathcal{X}_i^{(+;d)} = \{(p, I_{i,p})\}_{p \in \mathcal{P}^{(+)}(d)}$, $\mathcal{X}_i^{(-;d)} = \{(p, I_{i,p})\}_{p \in \mathcal{P}^{(-)}(d)}$.

View-Disease swap. Pick two fetal cases $i \neq j$ with $d_i \in \mathcal{D}_i$, $d_j \in \mathcal{D}_j$, and $d_i \neq d_j$. We swap only the salient views aligned by the established view-image correspondence: $(p, I_{i,p})$:

$$\tilde{\mathcal{X}}_{i \leftarrow j}^{(d_j)} \triangleq \mathcal{X}_i^{(-;d_j)} \cup \mathcal{X}_j^{(+;d_j)} \Big|_{\text{aligned by } (p, I_{i,p})}, \tilde{\mathcal{X}}_{j \leftarrow i}^{(d_i)} \triangleq \mathcal{X}_j^{(-;d_i)} \cup \mathcal{X}_i^{(+;d_i)} \Big|_{\text{aligned by } (p, I_{j,p})}. \quad (1)$$

Let x_i^{swap} and x_j^{swap} denote the full inputs (images + prompt) built from equation 1, i.e., $x_i^{\text{swap}} \triangleq (\tilde{\mathcal{X}}_{i \leftarrow j}^{(d_j)}, \text{prompt})$ and $x_j^{\text{swap}} \triangleq (\tilde{\mathcal{X}}_{j \leftarrow i}^{(d_i)}, \text{prompt})$. Note that any change in the images during swapping requires a synchronized update of the prompt accordingly. Our goal is to *reject* the receiver’s original disease set under swapped evidence. For each swapped input we form preference triplets: $(x_i^{\text{swap}}, \mathcal{D}_j, \mathcal{D}_i)$ and $(x_j^{\text{swap}}, \mathcal{D}_i, \mathcal{D}_j)$. The *chose* labels come from the donor and *reject* labels come from the receiver’s labels. We collect all triplets into the swap-derived set $\mathcal{D}_{\text{swap}}$. Early and mid-to-late pregnancy stages are swapped independently to account for their morphological differences.

Table 1: Comparison of **FetalMind** with other MLLM and unified multi-modal models on medical visual comprehension tasks. **Bold** and underlined text indicates the best performance and second-best performance, respectively. Note that * indicates models fine-tuned with *Supervised Fine-Tuning* to ensure a fair comparison.

Type	Model	#Params	Medical LVLM	NLG Metrics \uparrow				CE Metrics \uparrow			ACC \uparrow	Body F1-20 \uparrow	Avg. \uparrow
				B-1	B-4	MTR	R-L	P	R	F1			
w/o US Train	InternVL3	1B	\times	13.5	2.6	2.3	7.4	0.0	0.0	0.0	46.2	0.0	8.9
	QwenVL2.5	7B	\times	7.8	1.4	1.2	3.9	13.0	0.5	1.0	46.8	2.5	8.7
w/ US Train	BiomedGPT	182M	\checkmark	1.6	0.3	0.7	1.2	3.5	1.6	1.9	46.8	5.9	6.9
	LLaVA-Med	7B	\checkmark	0.9	0.3	0.4	0.6	2.0	0.1	0.2	46.2	0.8	5.6
	LLaVA-Med *	7B	\checkmark	6.3	3.0	4.4	5.6	1.9	0.1	0.1	46.9	0.8	11.6
	Med-Flamingo	8.3B	\checkmark	21.6	8.9	8.5	7.7	3.8	1.1	1.7	44.1	1.6	14.5
	Gemini 2.5 Pro	-	\times	16.9	7.0	9.9	12.9	19.4	16.1	17.6	71.4	26.4	24.2
	GPT-5	-	\times	28.3	8.3	4.8	12.4	19.1	12.6	15.2	71.6	23.6	24.1
	InternVL3 *	1B	\checkmark	14.1	4.0	4.9	6.5	<u>26.2</u>	18.9	22.0	78.2	39.9	23.9
	FetalMind-S1	1B	\checkmark	30.3	<u>9.2</u>	<u>15.5</u>	<u>12.4</u>	23.1	29.2	<u>25.8</u>	<u>79.0</u>	<u>45.2</u>	<u>29.7</u>
	FetalMind-M7	7B	\checkmark	33.9	23.1	30.4	30.7	34.7	<u>28.2</u>	31.1	81.3	50.2	38.2

Data-Centric Learning via SVPO. We optimize preference alignment on $\mathcal{D}_{\text{swap}}$ using Salient View Preference Optimization (SVPO). **The key idea is a strategy that builds preference pairs by mining Salient Views from knowledge graph on top of existing preference-optimization algorithms.** Either online rewards (e.g., PPO (Schulman et al., 2017)) or offline chosen/rejected pairs (e.g., DPO (Rafailov et al., 2024), CPO (Xu et al., 2024)) can be used; following prior visual alignment work (Yu et al., 2024a;b), we adopt the offline formulation. The SVPO objective is

$$\mathcal{L}_{\text{SVPO}}(\pi_{\theta}) = -\mathbb{E}_{(x, \mathcal{D}_w, \mathcal{D}_l) \sim \mathcal{X}} [\log \sigma(\beta(\log \pi_{\theta}(\mathcal{D}_w | x) - \log \pi_{\theta}(\mathcal{D}_l | x)))], \quad (2)$$

where σ is the sigmoid and $\beta > 0$ is a temperature. Let the contrastive score be $g \triangleq \log \pi_{\theta}(\mathcal{D}_w | x) - \log \pi_{\theta}(\mathcal{D}_l | x)$, $\Delta = \beta g$. The gradients are

$$\frac{\partial \mathcal{L}_{\text{prefer}}}{\partial \Delta} = \sigma(\Delta) - 1, \quad \frac{\partial \mathcal{L}_{\text{prefer}}}{\partial g} = \beta(\sigma(\Delta) - 1). \quad (3)$$

When the chosen and rejected responses are very close ($\Delta \approx 0$, i.e., hard pairs), $\sigma(\Delta) \approx \frac{1}{2}$ and hence $\frac{\partial \mathcal{L}_{\text{prefer}}}{\partial g} \approx -\frac{\beta}{2}$, providing a non-negligible signal that *simultaneously* increases $\log \pi_{\theta}(\mathcal{D}_w | x)$ and decreases $\log \pi_{\theta}(\mathcal{D}_l | x)$. Consequently, SVPO naturally emphasizes hard pairs and sharpens fine-grained distinctions (e.g., negation, units, laterality, anatomical loci) that are critical for medical report generation and diagnosis. As shown in Equation (2), SVPO reinforcement learning operates by constructing inputs x and pairing them with chosen samples \mathcal{D}_w and rejected samples \mathcal{D}_l . In our formulation, the training distribution is instantiated by the swap-derived dataset $\mathcal{D}_{\text{swap}}$.

Principles of Swap Construction. As shown in Figure 3c, we summarize four swap recipes for constructing preference pairs while preserving anatomical plausibility and inter-view consistency:

- ❶ **Disease-to-Normal.** Randomly sample two fetuses. For the receiver, remove disease-related images and replace them with the donor’s *normal* images for the corresponding views.
- ❷ **Normal-to-Disease.** Sample a normal receiver and an abnormal donor. Replace the receiver’s corresponding images with the donor’s *disease-related* images; if a corresponding plane is missing, append the donor’s disease-related plane set.
- ❸ **Disease-to-Disease.** Sample two abnormal fetuses with different disease. Remove the receiver’s disease-related images and insert the donor’s disease-related images to form a contrasted disease composition.
- ❹ **Disease Aggregation.** Sample two fetuses whose disease-related image sets are disjoint and merge them to synthesize a multi-disease case.

Global constraints. (1) Non-overlapping images are kept from the receiver rather than hallucinated. (2) When the number of images changes during a swap, the prompt must be updated accordingly.

5 EXPERIMENT

5.1 EXPERIMENTAL SETUP

Benchmarks. We randomly split data from nine centers into training/validation/test sets with a 7:1:2 ratio, and used data from the other three centers for external validation. To enable diverse

Table 2: Comparison of **FetalMind** with other LVLMs and unified multi-modal models on medical visual comprehension tasks. **Bold** and underlined indicates the best and second-best performance, respectively.

Model	Early Preg. \uparrow				Mid Preg. \uparrow				Late Preg. \uparrow			
	Micro-D	Macro-D	Micro-V	Macro-V	Micro-D	Macro-D	Micro-V	Macro-V	Micro-D	Macro-D	Micro-D	Macro-V
InternVL3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
QwenVL2.5	2.5	1.4	2.7	1.4	0.8	0.4	2.2	0.9	3.0	1.6	2.9	1.2
BiomedGPT	8.3	1.4	4.1	0.8	4.9	0.6	6.8	2.5	7.1	1.0	5.2	0.9
LLaVA-Med	0.9	0.2	0.8	0.2	0.5	0.1	1.2	0.2	0.3	0.1	0.0	0.0
LLaVA-Med *	0.4	0.1	0.4	0.1	0.1	0.0	1.1	0.2	0.7	0.1	0.5	0.1
Med-Flamingo	6.8	1.5	0.7	0.3	2.3	0.3	1.8	0.6	3.7	1.1	1.5	0.9
Gemini 2.5 Pro	20.5	13.8	21.4	19.6	19.5	16.2	27.2	17.2	24.5	17.7	27.4	16.5
GPT-5	13.4	6.9	14.1	12.5	17.9	14.8	25.7	18.3	21.3	14.9	24.1	17.2
InternVL3 *	25.1	19.6	37.2	11.1	23.2	7.9	41.3	13.2	24.1	15.6	38.7	11.1
FetalMind-S1	<u>25.8</u>	<u>30.7</u>	<u>27.8</u>	<u>18.5</u>	<u>30.2</u>	<u>19.3</u>	<u>47.9</u>	<u>21.6</u>	<u>36.9</u>	<u>30.2</u>	<u>44.5</u>	<u>18.1</u>
FetalMind-M7	41.0	36.0	44.5	20.3	35.2	22.1	52.1	22.0	39.6	31.6	49.6	18.7

evaluation, we extract gestational-age metadata from ultrasound reports and partition the test set into *early*, *mid*, and *late* subsets, assessing robustness and generalization across stages. The evaluation results confirm the performance improvements of our model, particularly evident in early pregnancy diagnosis and major malformations. The metrics are provided in Section D.

Baseline Methods. We compare **FetalMind** against nine MLLM baselines. InternVL3 (Zhu et al., 2025) and QwenVL-2.5 (Bai et al., 2025) were not trained on ultrasound data. The other seven models incorporate ultrasound in their training pipelines, including BiomedGPT (Zhang et al., 2024), LLaVA-Med (Li et al., 2023), Med-Flamingo (Moor et al., 2023b), Gemini 2.5 Pro, GPT-5, and our SFT variants LLaVA-Med* and InternVL3* fine-tuned on **FetalSigma-1M**. For open-source models, we evaluate the released checkpoints using their official prompting strategies. Although Gemini 2.5 Pro and GPT-5 do not explicitly disclose prenatal ultrasound data, their stable performance and reported medical pretraining suggest indirect exposure; we therefore categorize them as *with-ultrasound* in our analysis. Note that for models lacking native diagnostic capability, we obtain the corresponding diagnoses by passing their generated reports to GPT (Guo et al., 2025), using carefully crafted prompts together with a structured specification of the disease set.

Implementation Details. We train the model on NVIDIA A800 GPUs with one epoch for the alignment stage, three epochs for instruction tuning, and one epoch for reinforcement learning with SVPO. The learning rate is set to 5×10^{-5} , and the temperature parameter is fixed at $\beta = 0.0$. Our 1B model is instantiated from InternVL3, whereas the 7B variant is built upon Qwen2.5-VL. For fairness, we fix the image size to 224×224 for all models. More results are provided in Section A.

5.2 EVALUATION ON GENERAL MULTI-CENTER STUDY

Performance on Medical Diagnosis. Medical diagnosis requires accurate prediction of one or more standardized labels, directly impacting clinical decision-making and patient outcomes. On the twelve-center disease-classification benchmark (Table 1), **FetalMind-M7** improves binary abnormal/normal accuracy by 9.7%. Multi-label classification is particularly challenging for MLLMs because it demands disentangling subtle symptoms and mapping them to precise diagnoses. Under cross-entropy metrics, **FetalMind-M7** achieves an F1 gain of 13.5% and a recall gain of 9.3% over prior models. To further assess localization fidelity from diseases to fetal anatomy, we construct a disease-view mapping spanning 20 anatomical categories (e.g., cardiac, cerebral). As shown in the penultimate column of Table 1, **FetalMind** achieves a 23.8% gain, demonstrating the effectiveness of SED in grounding diseases to the correct images and reinforcing disease-view alignment.

Performance on Medical Report Generation. Medical report generation requires the model to generate a detailed report based on the provided medical scan. As shown in Table 1, **FetalMind-M7** achieves the best scores, outperforming strong baselines (e.g., Gemini 2.5 Pro and GPT-5) by approximately **+5.6%** (BLEU-1), **+14.2%** (BLEU-4), **+20.5%** (METEOR), and **+17.8%** (ROUGE-L). The lighter **FetalMind-S1** variant ranks second on most NLG metrics, indicating a favorable efficiency-performance trade-off. A visual comparison is provided in Figure 4. These gains suggest that SVPO encourages explicit correspondences between multiple images and diagnostic labels rather than treating images and labels as an undifferentiated set (see Figure 1b), thereby improving multi-image grounding and robustness for report generation and multi-label disease classification.

Performance on Different Stages of Pregnancy. Mastery of fetal ultrasound by physicians typically requires **3+ years** of education, considerably longer than X-ray interpretation (about 1 year),

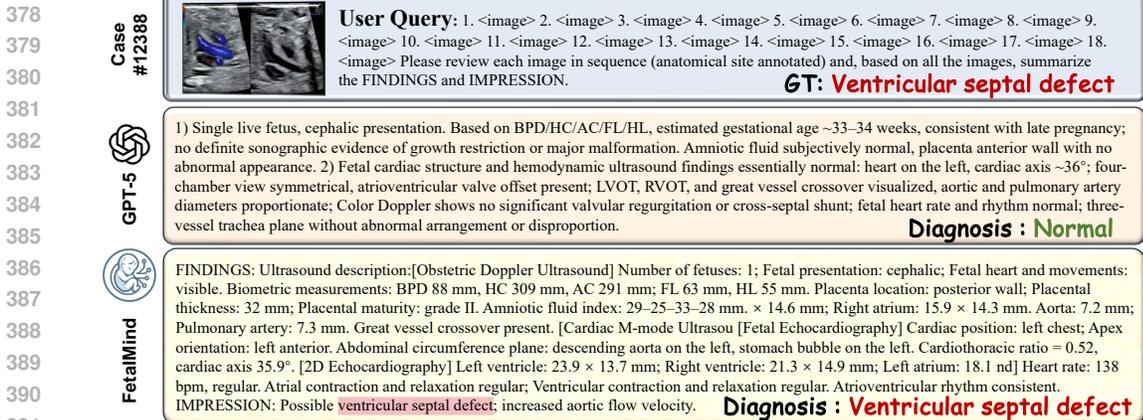


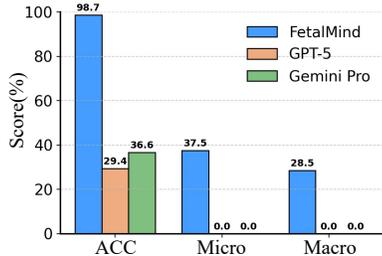
Figure 4: Illustration of *FetalMind* versus GPT-5 on a representative case (ID: 12388). The ground-truth diagnosis is a ventricular septal defect (VSD). GPT-5 misclassified the case as normal, likely due to its limited utilization of 2D and Doppler signals. In contrast, *FetalMind* correctly identified the VSD by integrating multi-view structural cues with blood-flow features. The report is truncated for brevity.

underscoring the task’s complexity. Following clinical practice, we stratify evaluation by gestational stage (*early, mid, late*) and report performance per trimester. As shown in Table 2, **Micro-D** denotes multi-label disease classification, while **Micro-V** measures performance after mapping diseases to anatomical regions. *FetalMind*-M7 surpasses all baselines across trimesters, with gains ranging from **2.2%** to **24.9%**, demonstrating strong generalization. Notably, in the *early* trimester, **Micro-D** improves by **20.5%**, highlighting the model’s value for earlier detection of fetal anomalies—enabling earlier, potentially actionable findings and affording more time for follow-up and clinical decision-making. [More experiments in Section A.6.](#)

5.3 EVALUATION ON THE NINE MAJOR MALFORMATIONS

To assess the model’s diagnostic capability for critical conditions, we curated 153 clinically confirmed cases covering nine major congenital anomalies, which are critical in prenatal ultrasound diagnosis in China, where misdiagnosis often leads to severe medical or legal consequences. These challenging cases were collected across three centers and multi-device models, providing clinically reliable ground-truth labels for evaluation. As shown in Figure 5, GPT-5 and Gemini 2.5 Pro, despite being state-of-the-art MLLMs for fetal ultrasound, consistently failed to identify these anomalies and often misclassified them as negative. In contrast, *FetalMind* achieved a diagnostic accuracy of **98%**, substantially surpassing all prior baselines across anomaly types and demonstrating robust decision support in complex clinical settings.

Figure 5: Diagnostic performance comparison in nine major malformations



5.4 ABLATION STUDIES

Ablation Studies on Strategy. As shown in Table 3, removing any of the three components with *token injection, spatial alignment, and SVPO* degrades performance. We summarize three key observations: **Obs. 1** Eliminating fetal token injection yields the smallest yet consistent drop across all metrics. This indicates that injecting fetal priors at the token level mainly strengthens fine-grained discrimination and stability, enabling the model to separate semantically similar but clinically distinct entities. **Obs. 2** Removing spatial alignment disproportionately reduces report generation quality while having a milder impact on diagnostic metrics. This suggests that cross-view spatial alignment primarily facilitates multi-image aggregation and narrative coherence, effectively multiple views into a *clinically interpretable* sum-

Table 3: Ablation study on *FetalMind* in the *FetalSigma-1M* dataset. The impact of without (w/o) and with (w) post-selection techniques.

Setting	B-4	F1	ACC	AVG
<i>FetalMind</i>	23.1	31.1	81.3	45.2
w/o Token inject	21.9	30.7	80.3	44.3
w/o Spatial align	16.3	29.4	80.6	42.1
w/o SED	13.7	26.7	80.1	40.5
w/ GRPO	9.7	24.2	79.2	37.3
w/ DPO	7.9	12.3	65.8	28.7
Vanilla	9.2	25.8	79.0	38.0

mary. **Obs. 6** Removing SED causes the largest overall decline, establishing it as the primary source of improvement. By aligning multi-view preferences, SED simultaneously enhances report readability and stabilizes diagnostic discrimination, underscoring its central role in multi-view reasoning.

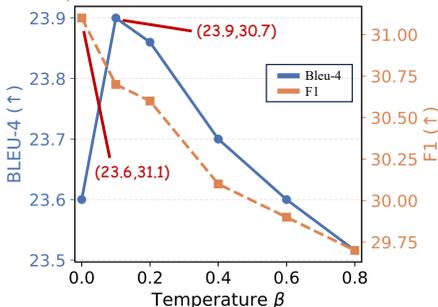
Ablation Studies on Reinforcement Learning. We further investigate the effect of different *reinforcement learning objectives* in Table 3. Compared with vanilla training, models optimized with DPO (Rafailov et al., 2023) or GRPO (Shao et al., 2024) perform worse across BLEU-4, F1, and ACC. In contrast, **FetalMind** achieves the strongest overall results. These findings underscore the importance of the post-selection procedure and demonstrate that SVPO with salient epistemic disentanglement is essential for enhancing diagnostic accuracy and producing clinically faithful reports.

5.5 PARAMETER SENSITIVITY ANALYSIS

Temperature β . As shown in Figure 6, we observe a distinct task-dependent trend. For *diagnostic classification*, lower temperatures consistently yield stronger performance, as reduced sampling stochasticity improves label consistency and raises F1/ACC. In contrast, for *report generation*, a mild degree of randomness proves beneficial: performance peaks around $\beta = 0.1$, balancing exploratory diversity with factual stability. These results suggest a near-deterministic setting for diagnosis and a small but nonzero temperature for narrative generation.

Report Generation vs. Diagnosis. **FetalMind** highlights an inherent heterogeneity between report generation and diagnostic classification in both task objectives and evaluation metrics. As shown in Figure 5, excessive determinism and insufficient randomness reduce report coverage and completeness. Enabling *controlled exploration* in lesion-related segments while preserving determinism for diagnostic-critical points, and adopting task-specific, temperature-aware inference, further improves overall performance.

Figure 6: Parameter sensitivity of temperature β in **FetalMind-M7**.



6 DISCUSSION

FetalMind achieves best performance on both fetal report generation and diagnostic, surpassing both general large models and domain-specific medical models. An insight emerges: structured tool usage in medical AI holds value. Compared with purely end-to-end methods, coupling the reasoning capacity of large models with domain basic modules consistently yields superior performance.

Generalists Versus Specialists. A notable finding is that general-purpose models (e.g., GPT-5, Gemini 2.5 Pro) overall outperform specialized medical models (e.g., LLaVA-Med (Li et al., 2024), Med-Flamingo (Moor et al., 2023b)). This indicates that narrow specialization may diminish the broad reasoning abilities conferred by large-scale pretraining. By integrating domain-specific tools under clinical guidance, **FetalMind** provides an effective bridge between the two paradigms.

Limitations & Future Work. Our evaluation remains retrospective and constrained by the available dataset, and prospective clinical studies are crucial for establishing real-world utility and safety. On the other hand, there remains a theoretical risk that the model may inadvertently learn “splicing artifacts” from synthetic data. Promising directions include: (1) tighter integration with PACS and ultrasound consoles for seamless clinical deployment; (2) uncertainty estimation and case triage to enhance clinician oversight; (3) broader coverage of rare anomalies and robustness to domain shift through active and continual learning; (4) privacy-preserving federated training across hospitals; and (5) extending disease-view graphs to temporal modalities. We anticipate that **FetalSigma-1M** and **FetalMind** will catalyze clinically grounded research toward trustworthy fetal ultrasound AI.

7 CONCLUSION

In this work, we present **FetalMind**, a clinically guided AI system for fetal ultrasound and, to our knowledge, the first unified framework addressing both report generation and diagnosis. By incorporating bipartite graph and disentangling disease-view heterogeneity, our SED aligns the model’s reasoning trajectory with real-world diagnostic workflows. Trained on the newly curated **FetalSigma-1M** comprising 20K reports from 12 centers, **FetalMind** consistently outperforms both open-source and proprietary baselines across all gestational stages. Beyond improvements, our findings underscore the critical role of structured clinical priors in building reliable AI systems.

486 REPRODUCIBILITY STATEMENT
487

488 To ensure the reproducibility of this research, we describe the experimental setup, data processing
489 steps, and key implementation details. Specifically, we employed reinforcement learning from the
490 MS-Swift framework and used LLaMA-Factory for supervised fine-tuning, with all implementations
491 developed in PyTorch. The datasets used in this work are derived from real clinical applications; a
492 subset of the reports is included in the paper, and we will also release the trained model weights.
493

494 REFERENCES
495

- 496 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
497 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
498 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 499 Rima Arnaout, Lara Curran, Yili Zhao, Jami C Levine, Erin Chinn, and Anita J Moon-Grady. An
500 ensemble of neural networks provides expert-level prenatal detection of complex congenital heart
501 disease. *Nature medicine*, 27(5):882–891, 2021.
- 502 Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani
503 Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. Openflamingo: An open-
504 source framework for training large autoregressive vision-language models. *arXiv preprint
505 arXiv:2308.01390*, 2023.
- 506 Reza Azad, Ehsan Khodapanah Aghdam, Amelie Rauland, Yiwei Jia, Atlas Haddadi Avval, Afshin
507 Bozorgpour, Sanaz Karimijafarbigloo, Joseph Paul Cohen, Ehsan Adeli, and Dorit Merhof. Med-
508 ical image segmentation review: The success of u-net. *IEEE Transactions on Pattern Analysis
509 and Machine Intelligence*, 2024.
- 510 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang,
511 Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan,
512 Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng,
513 Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv
514 preprint arXiv:2502.13923*, 2025.
- 515 Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones,
516 Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional AI:
517 Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- 518 JS Carvalho, R Axt-Fliedner, R Chaoui, JA Copel, BF Cuneo, D Goff, L Gordin Kopylov, K Hecher,
519 W Lee, AJ Moon-Grady, et al. Isuog practice guidelines (updated): fetal cardiac screening.
520 *Ultrasound Obstet Gynecol*, 61(6):788–803, 2023.
- 521 Ziming Cheng, Binrui Xu, Lisheng Gong, Zuhe Song, Tianshuo Zhou, Shiqi Zhong, Siyu Ren,
522 Mingxiang Chen, Xiangchao Meng, Yuxin Zhang, et al. Evaluating mllms with multimodal multi-
523 image reasoning benchmark. *arXiv preprint arXiv:2506.04280*, 2025.
- 524 Maria Chiara Fiorentino, Francesca Pia Villani, Mariachiara Di Cosmo, Emanuele Frontoni, and
525 Sara Moccia. A review on deep-learning algorithms for fetal ultrasound-image analysis. *Medical
526 image analysis*, 83:102629, 2023.
- 527 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
528 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
529 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 530 Yu Hou, Zaifu Zhan, and Rui Zhang. Benchmarking gpt-5 for biomedical natural language process-
531 ing. *arXiv preprint arXiv:2509.04462*, 2025.
- 532 Hongjie Hu, Hao Huang, Mohan Li, Xiaoxiang Gao, Lu Yin, Ruixiang Qi, Ray S Wu, Xiangjun
533 Chen, Yuxiang Ma, Keren Shi, et al. A wearable cardiac ultrasound imager. *Nature*, 613(7945):
534 667–675, 2023.
- 535 Thunakala Bala Krishna and Priyanka Kokil. Standard fetal ultrasound plane classification based on
536 stacked ensemble of deep learning models. *Expert Systems with Applications*, 238:122153, 2024.

- 540 Lok Hin Lee, Elizabeth Bradburn, Rachel Craik, Mohammad Yaqub, Shane A Norris, Leila Cheikh
541 Ismail, Eric O Ohuma, Fernando C Barros, Ann Lambert, Maria Carvalho, et al. Machine learning
542 for accurate estimation of fetal gestational age based on ultrasound images. *NPJ digital medicine*,
543 6(1):36, 2023.
- 544 Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Nau-
545 mann, Hoifung Poon, and Jianfeng Gao. Llava-med: Training a large language-and-vision assis-
546 tant for biomedicine in one day. *arXiv preprint arXiv:2306.00890*, 2023.
- 547 Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li.
548 LLaVA-NeXT-Interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv*
549 *preprint arXiv:2407.07895*, 2024.
- 550 Tianwei Lin, Wenqiao Zhang, Sijing Li, Yuqian Yuan, Binhe Yu, Haoyuan Li, Wanggui He, Hao
551 Jiang, Mengze Li, Xiaohui Song, et al. Healthgpt: A medical large vision-language model for
552 unifying comprehension and generation via heterogeneous knowledge adaptation. *arXiv preprint*
553 *arXiv:2502.09838*, 2025.
- 554 Ziyu Liu, Yuhang Zang, Xiaoyi Dong, Pan Zhang, Yuhang Cao, Haodong Duan, Conghui He, Yuan-
555 jun Xiong, Dahua Lin, and Jiaqi Wang. Mia-dpo: Multi-image augmented direct preference
556 optimization for large vision-language models. In *The Thirteenth International Conference on*
557 *Learning Representations*.
- 558 Fadillah Maani, Numan Saeed, Tausifa Saleem, Zaid Farooq, Hussain Alasmawi, Werner Diehl,
559 Ameera Mohammad, Gareth Waring, Saudabi Valappi, Leanne Bricker, et al. Fetalclip: A visual-
560 language foundation model for fetal ultrasound image analysis. *arXiv preprint arXiv:2502.14807*,
561 2025.
- 562 Michael Moor, Oishi Banerjee, Zahra Shakeri Hossein Abad, Harlan M Krumholz, Jure Leskovec,
563 Eric J Topol, and Pranav Rajpurkar. Foundation models for generalist medical artificial intelli-
564 gence. *Nature*, 616(7956):259–265, 2023a.
- 565 Michael Moor, Qian Huang, Shirley Wu, Michihiro Yasunaga, Cyril Zakka, Yash Dalmia, Ed-
566 uardo Pontes Reis, Pranav Rajpurkar, and Jure Leskovec. Med-flamingo: A multimodal
567 medical few-shot learner. July 2023b. URL <https://arxiv.org/abs/2307.15189>.
568 arXiv:2307.15189.
- 569 James P Neilson, Cochrane Pregnancy, and Childbirth Group. Ultrasound for fetal assessment in
570 early pregnancy. *Cochrane Database of Systematic Reviews*, 2010(1), 1996.
- 571 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong
572 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kel-
573 ton, Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan
574 Leike, and Ryan J. Lowe. Training language models to follow instructions with human feedback.
575 In *NeurIPS*, 2022.
- 576 John Pellerito, Bryann Bromley, Sandra Allison, Anil Chauhan, Stamatia Destounis, Eitan Dick-
577 man, Beth Kline-Fath, Joan Mastrobattista, Marsha Neumyer, Tatjana Rundek, et al. Aium-acr-
578 acog-smfm-sru practice parameter for the performance of standard diagnostic obstetric ultrasound
579 examinations. *Journal of Ultrasound in Medicine*, 37(11):E13–E24, 2018.
- 580 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
581 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
582 models from natural language supervision. In *International conference on machine learning*, pp.
583 8748–8763. PmLR, 2021.
- 584 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
585 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*
586 *in neural information processing systems*, 36:53728–53741, 2023.
- 587 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
588 Finn. Direct preference optimization: Your language model is secretly a reward model. In
589 *NeurIPS*, 2024.

- 594 LJ Salomon, Z Alfrevic, V Berghella, CM Bilardo, GE Chalouhi, F Da Silva Costa, E Hernandez-
595 Andrade, G Malinger, H Munoz, D Paladini, et al. Isuog practice guidelines (updated): perfor-
596 mance of the routine mid-trimester fetal ultrasound scan. *Ultrasound in Obstetrics and Gynecol-*
597 *ogy*, 59(6):840–856, 2022.
- 598 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
599 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 600
- 601 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
602 Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathemati-
603 cal reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 604
- 605 Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Mohamed Amin, Le Hou,
606 Kevin Clark, Stephen R Pfohl, Heather Cole-Lewis, et al. Toward expert-level medical question
607 answering with large language models. *Nature Medicine*, 31(3):943–950, 2025.
- 608 Saad Slimani, Salaheddine Hounka, Abdelhak Mahmoudi, Taha Rehad, Dalal Laoudiyi, Hanane
609 Saadi, Amal Bouziyane, Amine Lamrissi, Mohamed Jalal, Said Bouhya, et al. Fetal biometry and
610 amniotic fluid volume assessment end-to-end automation using deep learning. *Nature Communi-*
611 *cations*, 14(1):7047, 2023.
- 612 Lehan Wang, Haonan Wang, Honglong Yang, Jiayi Mao, Zehong Yang, Jun Shen, and Xiaomeng
613 Li. Interpretable bilingual multimodal large language model for diverse biomedical tasks. *arXiv*
614 *preprint arXiv:2410.18387*, 2024.
- 615
- 616 Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei Chen, Zhuang Liu, In So Kweon, and
617 Saining Xie. Convnext v2: Co-designing and scaling convnets with masked autoencoders. In
618 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16133–
619 16142, 2023.
- 620 Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton
621 Murray, and Young Jin Kim. Contrastive preference optimization: Pushing the boundaries of llm
622 performance in machine translation. *arXiv preprint arXiv:2401.08417*, 2024.
- 623
- 624 Tianyu Yu, Yuan Yao, Haoye Zhang, Taiwen He, Yifeng Han, Ganqu Cui, Jinyi Hu, Zhiyuan Liu,
625 Hai-Tao Zheng, Maosong Sun, et al. RIHF-V: Towards trustworthy mllms via behavior alignment
626 from fine-grained correctional human feedback. In *CVPR*, 2024a.
- 627 Tianyu Yu, Haoye Zhang, Yuan Yao, Yunkai Dang, Da Chen, Xiaoman Lu, Ganqu Cui, Taiwen He,
628 Zhiyuan Liu, Tat-Seng Chua, et al. RLAIIF-V: Aligning mllms through open-source ai feedback
629 for super gpt-4v trustworthiness. *arXiv preprint arXiv:2405.17220*, 2024b.
- 630 Kai Zhang, Rong Zhou, Eashan Adhikarla, Zhiling Yan, Yixin Liu, Jun Yu, Zhengliang Liu, Xun
631 Chen, Brian D Davison, Hui Ren, et al. A generalist vision–language foundation model for diverse
632 biomedical tasks. *Nature Medicine*, pp. 1–13, 2024.
- 633
- 634 Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen
635 Duan, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for
636 open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025.
- 637
- 638
- 639
- 640
- 641
- 642
- 643
- 644
- 645
- 646
- 647

APPENDIX

In this appendix, we provide supplementary material to further elucidate our approach. Section A expands on the experiments with detailed protocols and ablation studies. Section B introduces the preliminaries of the Salient Epistemic Disentanglement (SED) reinforcement learning module. Section C visualizes the standardized structured report template that guides fetal ultrasound report generation and diagnosis. Finally, Section D consolidates the evaluation metrics and their definitions used throughout the paper.

A MORE EXPERIMENTS

A.1 ATTENTION ANALYSIS

Implementation Details. We curate a total of 10,000 SVPO samples, with approximately 2,500 assigned to each of the four states. To mitigate confounding due to inter-institution variability, SED construction is restricted to within-center data. This choice is motivated by two practical considerations: ① report templates vary substantially across medical centers, introducing formatting and phrasing biases; and ② for a given fetus, all images are acquired on the same device at the same site. Constraining SED to a single center therefore attenuates center/device effects and yields a cleaner evaluation of SVPO behavior. **FetalMind was trained using data from multiple devices, including 15 types of ultrasound machines from over four manufacturers.**

To evaluate whether our proposed SED module indeed guides the model to focus on pathological regions after training, we conducted a quantitative attention analysis. Following the design in Figure 2 Left, we computed the MeanALLQ, defined as the mean attention weight over all query tokens across layers and heads, for both abnormal and normal ultrasound images. We then examined how often the attention allocated to abnormal images dominates that of normal images, thereby reflecting the model’s capacity to capture clinically salient cues. As summarized in Table 4, the baseline Qwen2.5-VL model achieves a dominance ratio of only 39.1% (713/1824). Incorporating additional training signals (Qwen2.5-VL*) improves this ratio to 52.4% (956/1824). In contrast, our FetalMind-M7 substantially outperforms both baselines, with abnormal images receiving higher attention weights in 80.7% of cases (1472/1824). These results clearly indicate that SED effectively enhances the model’s ability to attend to pathological regions, thus strengthening its diagnostic reliability.

Table 4: Ratio-based evaluation of attention dominance on salient images. The *Salient* denotes the number of abnormal cases with higher MeanALLQ values than normal cases, while the *Normal* is the total number of test cases. Percentages reflect the proportion of salient images receiving stronger attention. (*) indicates models further tuned with supervised fine-tuning (SFT).

Model	Salient	Normal	Percentage
Qwen2.5-VL-7B	713	1824	39.1%
Qwen2.5-VL-7B *	956	1824	52.4%
FetalMind-M7	1472	1824	80.7%

A.2 CONFUSION MATRIX

To further investigate the robustness of our framework and the fidelity of generated reports, we conducted additional retrospective evaluations involving clinical experts. Specifically, we compared two strong vision–language baselines, **Gemini 2.5 Pro** and **GPT-5**, alongside our method, to examine whether evaluators could distinguish model-generated reports from authentic clinical reports.

Figure 7 presents the aggregated confusion matrix across all 12 medical centers. Notably, evaluators often misclassified reports generated by large models as authentic, indicating that both Gemini 2.5 Pro and GPT-5 achieved a high level of realism in language style and clinical adequacy. Nevertheless, GPT-5 exhibited slightly higher indistinguishability, suggesting stronger alignment with clinical reporting conventions.

To further assess robustness under physiological heterogeneity, we stratified the evaluation by gestational stages. As illustrated in Figure 8, evaluator performance remained consistent across early-,

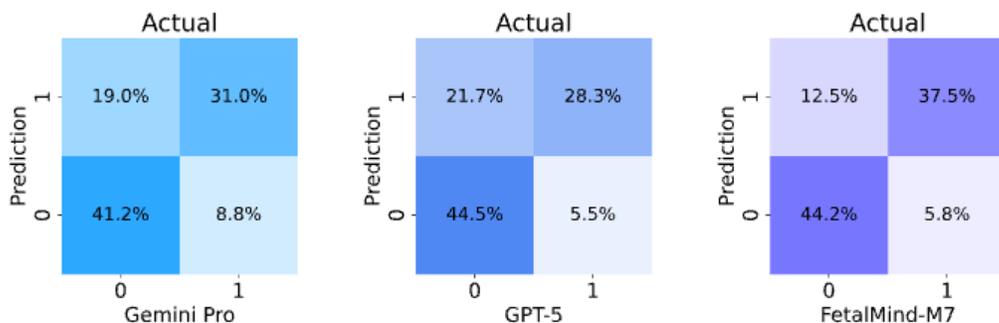


Figure 7: Confusion matrix for evaluators to identify reports generated by large models in the retrospective study, covering results from all 12 medical centers.

mid-, and late-gestation groups. The relative advantage of GPT-5 over Gemini 2.5 Pro persisted across all stages, reinforcing the conclusion that larger-scale alignment contributes to improved cross-condition fidelity. These findings collectively support the reliability of our framework and highlight the competitive performance of cutting-edge foundation models when benchmarked under rigorous human evaluation.

A.3 REPORT GENERATION STUDY

To further substantiate the effectiveness of our approach, we include a representative case study in Figure 9. In this example (Case #127858), the ground-truth diagnosis is *skeletal dysplasia*. While GPT-5 misclassifies the case as normal, **FetalMind** correctly identifies the pathology by jointly exploiting multi-view anatomical context and Doppler flow cues. This case illustrates how injecting domain-specific priors and explicitly modeling cross-view correspondences enables the system to recover subtle abnormalities that general-purpose LVLMs often overlook, thereby improving diagnostic reliability in fetal ultrasound.

A.4 GESTATIONAL AGE DISTRIBUTION

In addition to evaluator-based assessments, we also analyzed the distribution of gestational ages across centers in **FetalSigma-1M**. This is important because fetal ultrasound exhibits substantial heterogeneity in image appearance and reporting style at different stages of pregnancy, which may confound both training and evaluation if not carefully accounted for. Figure 10 shows the gestational age distributions extracted from three representative medical centers. Clear differences in case composition can be observed: while one center contributes a larger proportion of early-gestation cases, others are skewed toward mid-to-late gestation. Such heterogeneity motivates our stage-wise stratification strategy and provides empirical justification for evaluating model robustness under diverse physiological regimes. These analyses further highlight the challenges of building foundation models for fetal ultrasound and underline the necessity of multi-center, stage-aware evaluation.

A.5 REPORT CLASSIFICATION

Fetal Ultrasound Report Classification To validate the effectiveness of **FetalMind**, we introduce an ablation experiment where the model classifies fetal ultrasound reports based on a list of predefined disease labels. The process begins with the model generating a report from the ultrasound data, followed by selecting relevant disease labels based on the report’s content. The selected labels are then compared to the ground truth labels provided by clinical experts. The final classification accuracy is used to assess the model’s performance across several benchmarks. Our findings indicate that **FetalMind** offers a significant improvement in both diagnostic accuracy and clinical relevance compared to previous approaches. The prompt used to guide the model in classifying the ultrasound report is as follows:

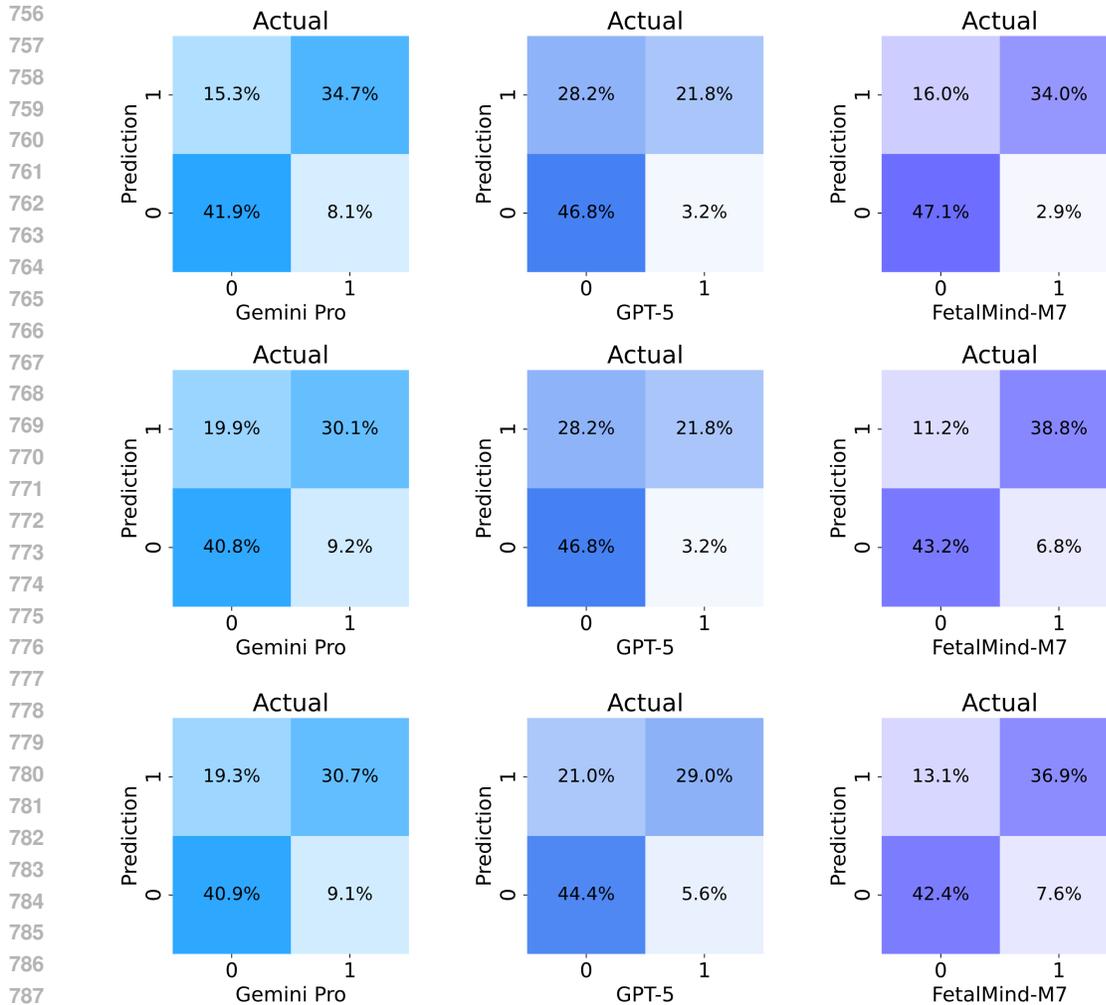


Figure 8: Confusion matrices illustrating evaluator performance in distinguishing reports generated by large language models during the retrospective study. Results are stratified by early-, mid-, and late-gestation stages, reflecting variability across different phases of pregnancy and highlighting the consistency of evaluation outcomes under diverse clinical conditions.

You are an expert in fetal ultrasound diagnosis. Based on the following ultrasound report, please select the disease labels that are explicitly mentioned or can be definitively inferred. The disease labels are provided in a predefined list.

The specific requirements are as follows:

1. Only select labels that are directly related to the content of the report.
2. If there are multiple disease labels, separate them with commas.
3. The output should be formatted as: Disease1, Disease2, ... (do not include numbering, explanations, or quotation marks).
4. If no disease labels are relevant, return an empty string.

Please review the report and select the disease labels accordingly.

Available Disease Labels: {Label1, Label2, Label3, ...}

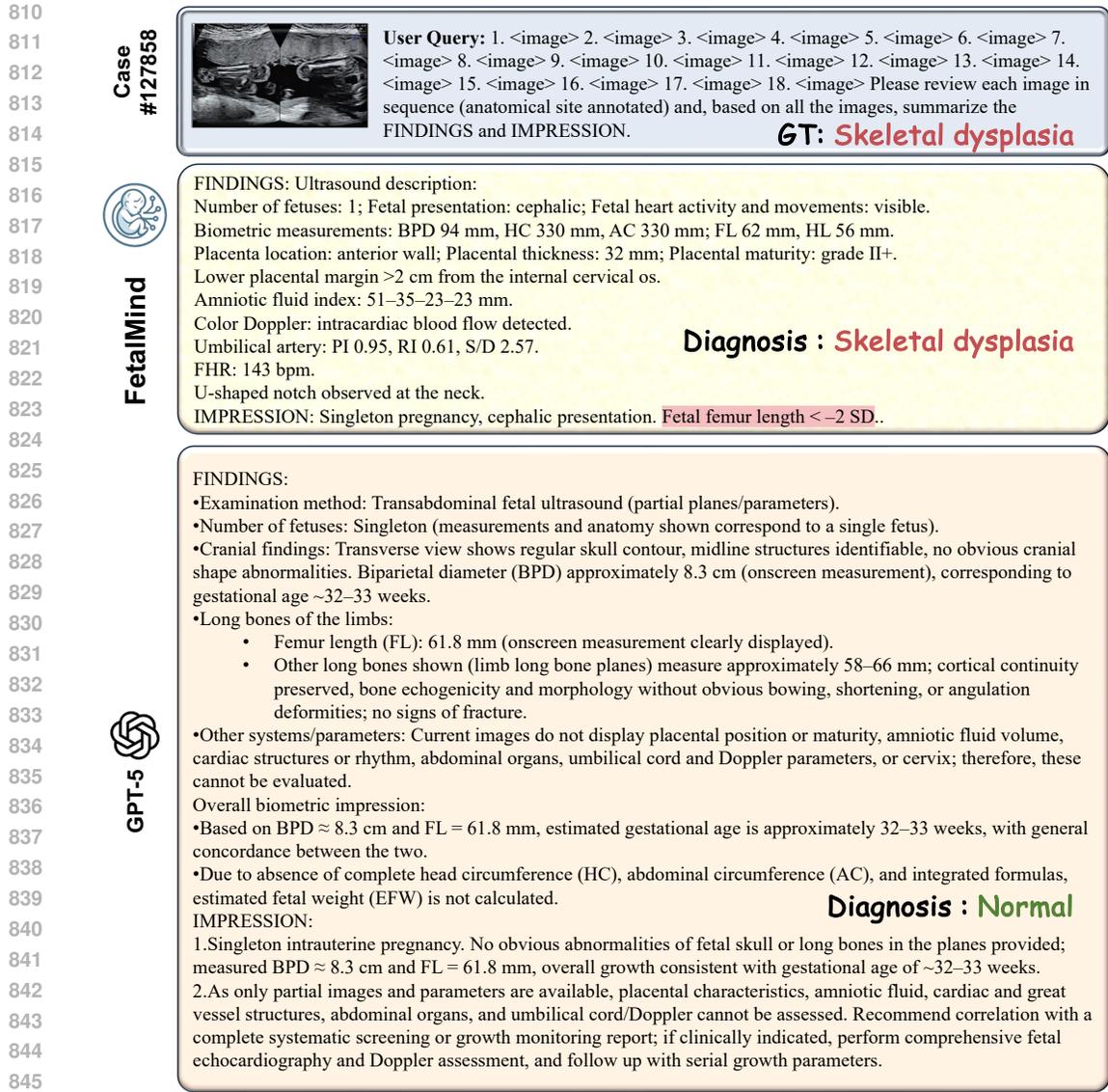


Figure 9: Illustration of FetalMind and GPT-5 Case Study. (Case 127858) Correct answer is skeletal dysplasia. GPT-5 misclassified it as normal, while FetalMind correctly identified skeletal dysplasia by integrating multi-view structures and blood flow features.

Table 5: Overall comparison of NLP and classification metrics between Doctor and Doctor+AI.

Method	BLEU-1	BLEU-4	ROUGE-1	ROUGE-L	METEOR	Precision _{micro}	F1 _{micro}
Doctor	75.388	67.817	77.450	72.019	40.592	0.568	0.562
Doctor+AI	88.532	81.605	86.002	85.717	59.351	0.679	0.653

Ultrasound Report : { [Insert ultrasound report here] }

Please provide the disease labels in the format mentioned above.

A.6 REAL-WORLD CLINICAL DECISION-MAKING ANALYSIS

To further validate the effectiveness of our method, we conducted a real-world clinical scenarios test on 56 cases from two centers (as shown in Table 5). Specifically, we set up three groups: one with a

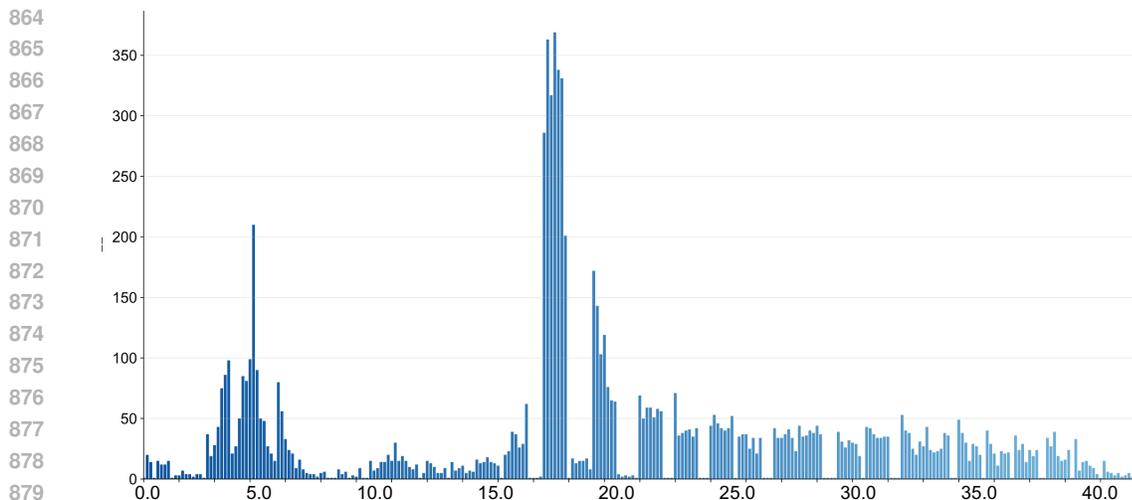


Figure 10: Visualization of gestational age distributions extracted from three medical centers. The figure highlights differences in case composition across centers, providing insights into data heterogeneity and supporting stratified analyses in subsequent model training and evaluation.

moderately experienced doctor, one with a moderately experienced doctor + AI, and a real clinical control group (three doctors including at least one highly experienced doctor). After completing the image collection for the examination, diagnosis was performed simultaneously, as shown in Table 1. As can be seen, our **FetalMind** can assist doctors by improving the neatness of the report writing and enhancing diagnostic accuracy.

A.7 VISUALIZATION OF THE DISEASE–VIEW GRAPH

In Figure 11, we present a visualization of the Disease–View graph using a Sankey diagram. This method effectively represents the relationships between different diseases and body regions, with the flow width indicating the intensity or frequency of each connection. Our disease–view bipartite graph contains 326 disease nodes, 54 view nodes, and 879 corresponding edges. All nodes are determined based on textbooks, clinical guidelines, and expert consensus, and subsequently standardized through unified terminology to ensure consistency. We further detail the expert-in-the-loop construction process: three clinicians with over 10 years of experience reviewed the preliminary disease–view relations, refined them, and conducted multiple rounds of discussion. For cases where expert opinions diverged, we clarify in the revised manuscript that the resolution followed a Delphi-style anonymous voting procedure or arbitration by a senior third expert.

B PRELIMINARY AND ANALYSIS

B.1 PRELIMINARY

To improve an LVLMM’s reasoning over *multi-image* inputs, we adopt *visual preference alignment*. This section formalizes the objective and uses *CPO* as a representative instantiation.

Visual Preference Alignment Preference alignment trains a model so that its output preferences conform to human (or proxy) preferences. Prominent paradigms include **Reinforcement Learning from Human Feedback (RLHF)** (Ouyang et al., 2022) and **Reinforcement Learning from AI Feedback (RLAIF)** (Bai et al., 2022). Let a dataset D consist of triplets $\{x, y_w, y_l\}$,¹ where x is a multimodal prompt—an interleaved sequence of images v and texts t —and y_w/y_l denote the *chosen* and *rejected* responses, respectively. Given a policy $\pi_\theta(y | x)$ and a reward model $r(x, y)$ that assigns higher scores to preferred responses, the visual preference alignment objective maximizes

¹For clarity we present single-sample notation; the extension to mini-batches is straightforward.

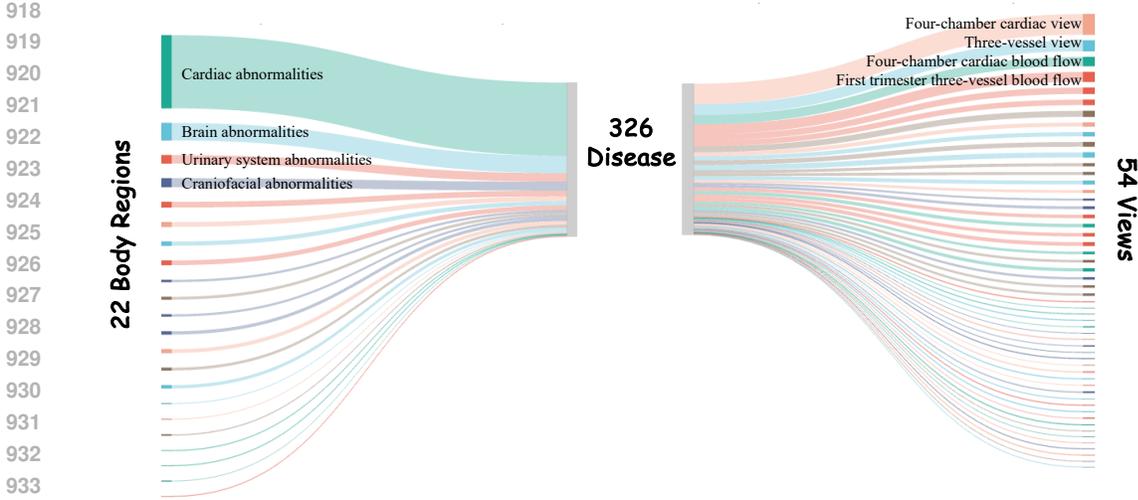


Figure 11: Visualization of Disease-View bipartite graph using a Sankey diagram. **Body Regions** represent different parts of the fetus, including the head, heart, and others.

expected reward:

$$\max_{\theta} \mathbb{E}_{x \sim D, y \sim \pi_{\theta}(y|x)} [r(x, y)], \quad (4)$$

where θ parameterizes the LVLMM. To mitigate overfitting and constrain drift from a reference policy π_{ref} , one augments the objective with a KL regularizer:

$$\max_{\theta} \left[\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(y|x)} [r(x, y)] - \beta D_{\text{KL}}(\pi_{\theta}(y | x) \| \pi_{\text{ref}}(y | x)) \right], \quad (5)$$

where $\beta > 0$ balances reward maximization and policy proximity. In practice, π_{ref} is the model snapshot before preference alignment.

CPO contrastive score CPO instantiates preference learning via a contrastive margin between the chosen and rejected responses:

$$\Delta = \beta(\log \pi_{\theta}(y_w | x) - \log \pi_{\theta}(y_l | x)), \quad \mathcal{L}_{\text{prefer}} = -\log \sigma(\Delta), \quad (6)$$

where $\sigma(\cdot)$ is the logistic sigmoid and $\beta > 0$ acts as a temperature.

Near-tie behavior (hard pairs) Let $g \triangleq \log \pi_{\theta}(y_w | x) - \log \pi_{\theta}(y_l | x)$ so that $\Delta = \beta g$. The gradients are

$$\frac{\partial \mathcal{L}_{\text{prefer}}}{\partial \Delta} = \sigma(\Delta) - 1, \quad \frac{\partial \mathcal{L}_{\text{prefer}}}{\partial g} = \beta(\sigma(\Delta) - 1). \quad (7)$$

For *hard* pairs where the two responses are nearly tied ($\Delta \approx 0$), we have $\sigma(\Delta) \approx \frac{1}{2}$ and thus

$$\frac{\partial \mathcal{L}_{\text{prefer}}}{\partial g} \approx -\frac{\beta}{2}, \quad (8)$$

yielding a substantial, stable signal that simultaneously increases $\log \pi_{\theta}(y_w | x)$ and decreases $\log \pi_{\theta}(y_l | x)$. This property encourages fine-grained discrimination among near-synonymous or subtly different responses—e.g., negation, units, laterality, or anatomical loci—crucial for medical report generation and diagnosis from multi-view ultrasound.

Difference from DPO DPO also optimizes a margin, but it uses a *reference-adjusted* form

$$\tilde{\Delta} = \beta \left[\log \frac{\pi_{\theta}(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right], \quad (9)$$

which entangles the learning signal with the quality and stylistic biases of π_{ref} and typically incurs additional compute/memory overhead. In contrast, CPO’s margin depends solely on π_{θ} , delivering a cleaner, reference-free signal on near-ties and promoting a more compact, clinically faithful chosen-response distribution for multi-image inputs.

Table 6: Ablation study on `FetalMind` in the `FetalSigma-1M` dataset. The impact of without (w/o) and with (w) post-selection techniques.

Setting	B-4	F1	ACC	AVG
<code>FetalMind</code>	23.1	31.1	81.3	45.2
w/o SED	13.7	26.7	80.1	40.5
w/ GRPO	9.7	24.2	79.2	37.3
w/ DPO	7.9	12.3	65.8	28.7
Vanilla	9.2	25.8	79.0	38.0

B.2 ANALYSIS OF SVPO AND SED

SVPO is readily extensible to multimodal large language models (MLLMs) that perform multi-image reasoning, and can be seamlessly integrated into multi-image inference pipelines according to task requirements. By explicitly distinguishing salient from non-salient images, SVPO improves both computational efficiency and predictive accuracy. For example, in a three-image joint analysis scenario where the key evidence primarily resides in Image 1, SVPO effectively steers the model toward the most informative visual cues.

SED introduces graph-aware reasoning by leveraging a bipartite graph to separate salient from non-salient images. When combined with SVPO, SED further establishes preference relations between abnormal images and target conditions, allowing condition–image structural information to be naturally injected into the MLLM’s reasoning process. This design closely mirrors clinical workflows, where physicians select key views, focus on abnormal regions, and integrate disease knowledge across multiple images. Consequently, the framework is particularly well-suited to multi-image reasoning tasks with explicit graph-structured relationships.

In summary, SED embeds and strengthens SVPO, enabling the model not only to capture saliency relationships across images, but also to perform condition–image relational reasoning via graph structures. This yields a more principled foundation for interpretability and reliability in structured medical applications of MLLMs.

B.3 ANALYSIS OF REINFORCEMENT LEARNING METHODS

DPO performance drop causes and the role of CPO’s BC regularizer. DPO relies on a fixed reference policy π_{ref} (typically an SFT model) and optimizes a preference loss of the form $L(\pi_{\theta}; \pi_{\text{ref}})$. This implicitly constrains π_{θ} to stay close to the reference, which can be suboptimal in our setting: in specialized domains such as medicine, the reference model often underfits domain-specific knowledge, and hard-anchoring π_{θ} to such a reference can limit the achievable performance.

In contrast, CPO removes this dependence by setting $\pi_{\text{ref}} = U$, i.e., a uniform reference, and directly optimizing the contrastive objective $L(\pi_{\theta}, U)$. This design allows the policy to move beyond the limitations of the reference model and better align with the preference data. However, a purely contrastive preference loss $L(\pi_{\theta}, U)$ only encodes *relative* signals (“ y_w is preferred over y_l ”), without constraining the *absolute* likelihood of preferred outputs. As a result, optimizing only the preference term can drive the model toward over-emphasizing superficial or stylistic characteristics of the preferred responses, rather than preserving factual correctness and faithfulness. In other words, the model may learn that “more elaborate / more confident / more verbose” is preferred, and lean toward “stylistic enhancement” instead of robustly modeling the underlying target distribution.

To address this, we introduce a BC regularizer

$$\mathbb{E}_{(x, y_w) \sim D} KL(\pi_w(y_w|x) \parallel \pi_{\theta}(y_w|x)) < \epsilon,$$

which, as shown in Appendix C, is equivalent to adding a negative log-likelihood term on the preferred data:

$$L_{\text{CPO}}(\pi_{\theta}) = L(\pi_{\theta}, U) - \mathbb{E}_{(x, y_w) \sim D} [\log \pi_{\theta}(y_w|x)].$$

This BC term brings two concrete benefits:

- (i) **Preventing divergence from the true preferred data distribution.** The BC regularizer anchors π_{θ} to the empirical distribution of preferred samples, preventing the policy from drifting too far away from what is actually observed in the data. This mitigates the risk

of probability mass collapsing onto overly confident or stylistically extreme outputs and stabilizes training, especially in the absence of a strong reference policy.

- (ii) **Providing an absolute learning signal beyond relative comparisons.** While the preference loss $L(\pi_\theta, U)$ only tells the model that “ y_w is better than y_l ”, the BC term directly encourages high likelihood on y_w itself. This provides an *absolute* supervised signal on preferred outputs, complementing the purely contrastive objective and ensuring that the model learns not only which response is better, but also what a good response should look like in distributional terms.

These properties are particularly important in the medical domain, where the target behavior is highly deterministic and correctness is much more critical than stylistic variation. In such a setting, it is not sufficient to merely prefer one response over another; the model must consistently produce stable, factually accurate outputs that closely match expert-like references. The BC regularizer is therefore especially well-suited here, as it pulls the model toward a sharp, well-calibrated distribution over medically correct responses rather than encouraging diversity or style.

This interpretation is also consistent with our empirical analysis. As shown in Figure 6 of the main text, using a lower sampling temperature leads to better performance. This observation aligns with the role of the BC term: by encouraging higher likelihood on preferred responses, CPO effectively shapes a sharper and more deterministic output distribution, which is desirable in high-stakes medical applications. Together, these theoretical and empirical considerations justify our choice of CPO with BC regularization over standard DPO in the proposed framework.

GRPO Performance Degradation. In our fetal ultrasound experiments, we observe that GRPO-based reinforcement learning yields performance degradation compared to supervised models. The root cause is that the conventional rewards optimized by GRPO act only as imperfect proxies for real clinical objectives. Such proxy rewards fail to capture fine-grained anatomical consistency, multi-view joint reasoning, and standardized report structures. The policy therefore tends to overfit these incomplete signals—exploiting phrasing patterns, templates, or other superficial regularities to “game the reward”—while degrading on clinically critical attributes such as localization accuracy, measurement validity, and structural coverage.

Moreover, GRPO updates policies by sampling candidate responses and optimizing based on relative rewards, which introduces additional stochasticity and high-variance gradients. This perturbs the likelihood distribution learned by supervised training—a distribution that is already close to optimal for the near-deterministic mapping required in fetal ultrasound—and drives the policy toward a small set of “reward-seeking” modes, reducing robustness and generalization.

A promising direction is to train a dedicated reward model that evaluates each prediction and provides more clinically aligned feedback, supplying GRPO with learning signals that better reflect real diagnostic criteria. This approach is particularly compelling in complex fetal ultrasound settings that require multi-image reasoning and coverage across diverse anatomical structures and conditions.

B.4 ANALYSIS OF GPT-BASED VS. DIRECT DIAGNOSIS

For baseline models that do not provide native diagnostic outputs, we adopt a two-step evaluation protocol in which GPT is used to infer diagnostic labels from the generated reports. This indirect procedure can introduce additional noise, since inaccuracies in GPT’s second-step extraction may lead to an artificial underestimation of baseline performance. To quantify this effect, we apply the same two-step evaluation to `FetalMind` reports. As shown in the Table 7, the resulting performance degradation is within 1%. Here are several potential sources of error:

- **Misinterpretation by GPT.** GPT may misread the semantics of free-text reports, especially when negations or rhetorical expressions are involved. For instance, “atrial septal defect” can be misclassified as “ventricular septal defect”. Our Fetal Token Injection module is explicitly designed to mitigate this issue by introducing special tokens for key anatomical and pathological terms, thereby reducing semantic confusion at the tokenization level.
- **Ambiguous degree modifiers in medical language.** Clinical descriptions frequently include degree-related qualifiers such as “mild” or “suspicious”. These modifiers may be

Table 7: Analysis of GPT-Based vs. Direct Diagnosis. **Bold** and underlined text indicates the best performance and second-best performance, respectively. Note that * indicates models fine-tuned with *Supervised Fine-Tuning* to ensure a fair comparison.

Type	Model	#Params	GPT Diagnosis	CE Metrics \uparrow			ACC \uparrow	Body F1-20 \uparrow
				P	R	F1		
w/ US Train	Gemini 2.5 Pro	-	✓	19.4	16.1	17.6	71.4	26.4
	GPT-5	-	✓	19.1	12.6	15.2	71.6	23.6
	InternVL3 *	1B	✗	26.2	18.9	22.0	78.2	39.9
	FetalMind-S1	1B	✗	23.1	29.2	25.8	79.0	45.2
	FetalMind-M7	7B	✗	34.7	<u>28.2</u>	31.1	81.3	50.2
	FetalMind-M7	7B	✓	<u>34.2</u>	27.6	<u>30.7</u>	<u>80.6</u>	<u>50.0</u>

interpreted inconsistently by LLMs under different contexts, leading to over-calling or under-calling certain findings.

- **Heterogeneous reporting styles across centers.** The same underlying condition can be phrased very differently by radiologists at different institutions. For example, Center A might report “Choroid plexus cyst noted in left ventricle”, whereas Center B might describe “Anechoic lesion detected within the choroid plexus”. Although both correspond to a choroid plexus cyst, the lexical variation can introduce additional challenges for robust automatic parsing.

B.5 INVESTIGATION OF ERROR SAMPLES

As shown in Fig. Figure 12, we qualitatively analyze representative failure cases of our model on fetal ultrasound. We observe two main error patterns:

- **Over-sensitivity to minor findings.** In small or borderline examinations, the model can be overly sensitive to subtle variations, occasionally assigning pathological labels to findings that experienced clinicians would still consider within normal limits.
- **Inter-center inconsistency in annotation standards.** Different centers may adopt slightly different criteria for the same condition, for example using different thresholds to define increased nuchal translucency (e.g., > 2.5 mm vs. > 2.6 mm). This issue can be mitigated by harmonizing disease labels according to international guidelines and remapping site-specific criteria to a unified standard.

C TRAINING TEMPLATE

C.1 FETAL ULTRASOUND REPORT TEMPLATE

To promote both clinical validity and cross-center consistency, we constructed a standardized obstetric ultrasound report template by systematically consolidating and harmonizing recommendations from multiple international guidelines, including those issued by the ISUOG, AIUM, and Chinese Medical Association. As illustrated in Figure 13 and Figure 14, we release both an English and a Chinese version of the template. The English version facilitates alignment with widely adopted global standards, while the Chinese version ensures applicability in large-scale domestic clinical practice. Together, these templates provide a unified and clinically grounded structure for report writing, enabling reliable data annotation, model training, and evaluation. Importantly, by establishing a guideline-based framework, the templates mitigate variability across institutions and languages, offering a scalable foundation for developing deep learning systems that generalize robustly across centers, devices, and populations.

1134 C.2 INSTRUCTION CONTENT FOR CLINICAL EXPERT REFERENCE
1135

1136 Below, we present instruction templates for both report generation and diagnostic reasoning. These
1137 templates establish a consistent and structured reference framework for clinical experts during model
1138 evaluation, ensuring that model outputs are assessed according to unified and standardized criteria.

1139 **📄 Diagnosis Template (Expert Reference)**

1140 1. <Aortic arch view><image>2. <Four-chamber view with intracardiac echogenic
1141 focus><image>3. <Superior-inferior vena cava view><image>4. <Three-
1142 vessel view><image>5. <Four-chamber view><image>6. <Femur long-axis
1143 view><image>7. <Biparietal diameter plane><image>8. <Thoracic coronal
1144 view><image>. Please provide a professional fetal assessment based on the anatomical
1145 ultrasound planes and indicate whether one or more abnormalities are present.
1146
1147

1148 **📄 Report Template (Expert Reference)**

1149 1. <Biparietal diameter plane><image>2. <Longitudinal view of the spine><image>3.
1150 <Cerebellar plane><image>4. <Median sagittal view of the fetal face><image>5.
1151 <Lateral ventricle plane><image>6. <Transverse abdominal plane through the blad-
1152 der><image>7. <Coronal view of both kidneys><image>8. <Interorbital dist-
1153 ance plane><image>9. <Transverse view of both kidneys><image>10. <Sagittal
1154 view of a single kidney><image>11. <Transverse abdominal diameter measure-
1155 ment plane><image>12. <Three-vessel view><image>13. <Coronal view of the fin-
1156 gers><image>14. <Left ventricular outflow tract view><image>15. <Four-chamber
1157 heart view><image>16. <Pulmonary artery bifurcation><image>17. <Coronal
1158 view of the lips><image>18. <Gallbladder view><image>19. <Femur long-axis
1159 view><image>. Please carefully review the above ultrasound images (the image order
1160 and anatomical regions have been specified). After an overall assessment, please provide a
1161 detailed description of the Findings and Impression.
1162
1163
1164

1165 **D EVALUATION METRICS**
1166

1167 In this section, we provide a detailed mathematical formulation of common metrics used for eval-
1168 uating Natural Language Generation (NLG) tasks and Classification Evaluation (CE) tasks. These
1169 metrics, such as BLEU, METEOR, ROUGE-L, Precision, Recall, and F1-Score, are used to assess
1170 the quality and effectiveness of machine-generated text in comparison to ground truth references.
1171

1172 **D.1 BLEU (B-1 AND B-4)**
1173

1174 BLEU (Bilingual Evaluation Understudy) measures the precision of n-grams between the generated
1175 and reference texts. It is often used for machine translation and other NLG tasks. BLEU considers
1176 the precision of unigrams (B-1) and 4-grams (B-4), calculating the overlap between the generated
1177 text and reference texts.
1178

$$1179 \text{B-1} = \text{Precision}_1 = \frac{\sum_{n=1}^N \text{Count}_{\text{match},1}}{\sum_{n=1}^N \text{Count}_{\text{generated},1}} \quad (10)$$

$$1180 \text{B-4} = \text{Precision}_4 = \frac{\sum_{n=1}^N \text{Count}_{\text{match},4}}{\sum_{n=1}^N \text{Count}_{\text{generated},4}} \quad (11)$$

1181 Where: - $\text{Count}_{\text{match},n}$ represents the number of n-grams that appear in both the reference and the
1182 generated text. - $\text{Count}_{\text{generated},n}$ represents the total number of n-grams in the generated text.
1183
1184
1185
1186
1187

1188 BLEU can be extended with a brevity penalty (BP) to account for the length of the generated text:
1189

$$1190 \text{ BLEU} = \text{BP} \times \exp \left(\sum_{n=1}^N w_n \log p_n \right) \quad (12)$$

1191
1192 Where w_n is the weight for each n-gram (usually uniform), and p_n is the precision of n-grams of
1193 size n .

1194 D.2 METEOR (MTR)

1195 METEOR (Metric for Evaluation of Translation with Explicit ORdering) improves upon BLEU by
1196 incorporating synonymy, stemming, and word-order preservation. METEOR balances precision and
1197 recall with an F-score, considering the meaning of words (synonyms) and morphological variations
1198 (stemming).

$$1200 \text{ MTR} = \text{F}(\text{Precision, Recall, Synonymy, Stemming}) \quad (13)$$

1201 Where: - Precision is the proportion of generated words that match the reference words. - Recall is
1202 the proportion of reference words that match the generated words. - Synonymy adjusts for synonyms
1203 (i.e., different words with similar meanings). - Stemming adjusts for different forms of the same
1204 word (e.g., "running" vs. "run").

1205 The F-measure is used to combine precision and recall:

$$1206 \text{ F(P, R)} = \frac{10 \cdot \text{P} \cdot \text{R}}{9 \cdot \text{P} + \text{R}} \quad (14)$$

1210 D.3 ROUGE-L (R-L)

1211 ROUGE (Recall-Oriented Understudy for Gisting Evaluation) is a set of metrics primarily used for
1212 evaluating machine-generated summaries. The ROUGE-L metric focuses on the longest common
1213 subsequence (LCS) between the reference and generated text, which captures the order of the words.

1214 The ROUGE-L score is calculated as:

$$1215 \text{ R-L} = \frac{\text{LCS}(\text{generated, reference})}{\text{Length of reference}} \quad (15)$$

1216 Where $\text{LCS}(\text{generated, reference})$ is the length of the longest common subsequence between the
1217 generated text and the reference text. The LCS metric encourages the preservation of word order,
1218 which is crucial for the quality of text generation.

1219 Additionally, ROUGE can be extended to compute recall (R) and precision (P) as follows:

$$1220 \text{ R} = \frac{\text{LCS}}{\text{Length of reference}}, \quad \text{P} = \frac{\text{LCS}}{\text{Length of generated text}} \quad (16)$$

1230 D.4 PRECISION (P)

1231 Precision is a metric used in classification tasks, which measures the accuracy of the predictions by
1232 comparing the true positives (TP) to the total predicted positives (TP + FP):

$$1233 \text{ P} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (17)$$

1234 Where: - TP represents the number of true positive instances (correctly predicted relevant instances).
1235 - FP represents the number of false positive instances (incorrectly predicted relevant instances).

1242 D.5 RECALL (R)

1243
1244 Recall measures how well the classifier identifies all relevant instances by comparing the true posi-
1245 tives (TP) to the total number of actual positives (TP + FN):

$$1246 R = \frac{TP}{TP + FN} \quad (18)$$

1247
1248 Where: - FN represents the number of false negative instances (relevant instances that were incor-
1249 rectly predicted as irrelevant).

1252 D.6 F1 SCORE (F1)

1253
1254 The F1 Score is a harmonic mean of precision and recall, providing a balanced measure of classifi-
1255 cation performance. It is particularly useful when dealing with imbalanced datasets:

$$1256 F1 = 2 \times \frac{P \times R}{P + R} \quad (19)$$

1257
1258 The F1 Score is maximized when both precision and recall are high, making it an excellent metric
1259 when both false positives and false negatives are equally important.

1263 D.7 MACRO AND MICRO AVERAGING FOR PRECISION, RECALL, AND F1

1264
1265 In multi-class classification tasks, we often calculate macro and micro averages for precision, recall,
1266 and F1 score:

1267 **Macro Average:** The macro average treats all classes equally by averaging the individual scores of
1268 each class:

$$1269 \text{Macro } P = \frac{1}{C} \sum_{i=1}^C P_i, \quad \text{Macro } R = \frac{1}{C} \sum_{i=1}^C R_i, \quad \text{Macro } F1 = \frac{1}{C} \sum_{i=1}^C F1_i \quad (20)$$

1270
1271 Where C is the number of classes, and P_i , R_i , and $F1_i$ are the precision, recall, and F1 scores for
1272 class i .

1273
1274 **Micro Average:** The micro average aggregates the true positives, false positives, and false negatives
1275 across all classes and then calculates the precision, recall, and F1:

$$1276 \text{Micro } P = \frac{\sum_{i=1}^C TP_i}{\sum_{i=1}^C (TP_i + FP_i)}, \quad \text{Micro } R = \frac{\sum_{i=1}^C TP_i}{\sum_{i=1}^C (TP_i + FN_i)}, \quad (21)$$

$$1277 \text{Micro } F1 = 2 \times \frac{\text{Micro } P \times \text{Micro } R}{\text{Micro } P + \text{Micro } R}$$

1278
1279 Where TP_i , FP_i , and FN_i are the true positives, false positives, and false negatives for class i ,
1280 respectively.

1287 E THE USE OF LARGE LANGUAGE MODELS (LLMS)

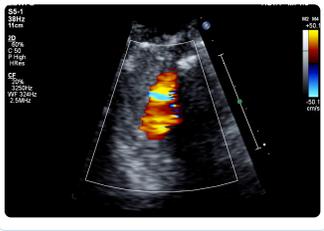
1288
1289 During manuscript preparation, we employed large language models (LLMs), specifically GPT-5,
1290 strictly as writing assistants to enhance grammar, clarity, and readability. Their role was limited to
1291 rephrasing for improved flow and correcting typographical errors. The scientific ideas, experimental
1292 design, analyses, and conclusions were conceived and developed entirely by the human authors. All
1293 model-generated text was carefully reviewed and edited by the authors, who take full responsibility
1294 for the manuscript’s accuracy and originality.

1296				
1297				
1298				
1299				
1300				
1301				
1302				
1303				
1304	Case #3329	User Query: 1. <image> 2. <image> 3. <image> 4. <image> 5. <image>. Please review each image in sequence (anatomical site annotated) and, based on all the images, summarize the FINDINGS and IMPRESSION.		
1305		Ground Truth	FINDINGS: Ultrasound description: Fetal heart and movements: visible; cephalic presentation. Anatomical structures: spine visualized.	
1306			FetalMind	FINDINGS: Ultrasound description: Fetal heart and movements: visible. Left renal pelvis: 7.5 mm; right renal pelvis: 10.1 mm.
1307				IMPRESSSION: Diagnosis : Normal
1308				IMPRESSSION: Diagnosis: Bilateral renal pelvic dilatation
1309				Spine re-examination: visualized.
1310	Bilateral renal pelvic dilatation.			
1311				
1312				
1313	Case #34212	User Query: <image> 2. <image> 3. <image> 4. <image> 5. <image> 6. <image> 7. <image> 8. <image> 9. <image> 10. <image> 11. <image> 12. <image> 13. <image> 14. <image> 15. <image> 16. <image>. Please review each image in sequence (anatomical site annotated) and, based on all the images, summarize the FINDINGS and IMPRESSION.		
1314		Ground Truth	FINDINGS: Ultrasound description: [Transabdominal] Number of fetuses: 1; Fetal position: cephalic; Fetal heart activity and movements: present. [Fetal biometric measurements] Crown-rump length (CRL): 73 mm; Biparietal diameter: 23 mm; Head circumference: 86 mm; Abdominal circumference: 65 mm; Femur length: 9 mm; Humerus length: 8 mm. Diagnosis : Normal	
1315			FetalMind	FINDINGS: Ultrasound description: [Transabdominal] Number of fetuses: 1; Fetal position: cephalic; Fetal heart activity and movements: present. [Fetal biometric measurements] Crown-rump length (CRL): 73 mm; Biparietal diameter: 23 mm; Head circumference: 86 mm; Abdominal circumference: 66 mm; Femur length: 8 mm; Humerus length: 8 mm. Diagnosis: Increased nuchal translucency
1316				Placental location: anterior wall and right lateral wall; placental thickness: 15 mm; placental maturity: grade 0; maximum amniotic fluid pocket depth: 32 mm. [Nuchal translucency (NT) thickness]: 2.6 mm. The following anatomical structures are visualized: Head: cranial halo; midline of the brain; "butterfly sign"; left orbit; right orbit. Chest: thoracic cavity. Abdomen: abdominal wall; gastric bubble; bladder; umbilical cord vessels: 3. Limbs: left femur; left tibia and fibula; right femur; right tibia and fibula; left humerus; left radius and ulna; left wrist; right humerus; right radius and ulna; right wrist. [M-mode ultrasound] Fetal heart rate: 157 beats per minute. [Uterine arteries] Uterus appears in a pregnant state. On the medial side of the right pelvic wall, the right uterine artery runs upward and forward, crossing the right external iliac artery; PI: 0.64; RI: 0.46; early diastolic notch: not seen. On the anteromedial side of the left pelvic wall, the left uterine artery runs upward and forward, crossing the left external iliac artery; PI: 1.82; RI: 0.77; early diastolic notch: not seen. IMPRESSION: Singleton, viable fetus. Fetal size corresponds to 13 weeks 3 days of gestation.
1317				Placental location: anterior wall; placental thickness: 13 mm; placental maturity: grade 0; maximum amniotic fluid pocket depth: 31 mm. [Nuchal translucency (NT) thickness]: 2.6 mm. The following anatomical structures are visualized: Head: cranial halo; midline of the brain; "butterfly sign"; left orbit; right orbit. Chest: thoracic cavity. Abdomen: abdominal wall; gastric bubble; bladder; umbilical cord vessels: 3. Limbs: left femur; left tibia and fibula; right femur; right tibia and fibula; left humerus; left radius and ulna; left wrist; right humerus; right radius and ulna; right wrist. [M-mode ultrasound] Fetal heart rate: 159 beats per minute. [Uterine arteries] Uterus appears in a pregnant state. On the medial side of the right pelvic wall, the right uterine artery runs upward and forward, crossing the right external iliac artery; PI: 1.31; RI: 0.69; early diastolic notch: not seen. On the anteromedial side of the left pelvic wall, the left uterine artery runs upward and forward, crossing the left external iliac artery; PI: 1.32; RI: 0.66; early diastolic notch: not seen. IMPRESSION: Singleton pregnancy. Fetal size corresponds to 13 weeks 3 days of gestation. Increased nuchal translucency (NT)
1318				
1319				
1320				
1321				
1322				
1323				
1324				
1325				
1326				
1327				
1328				
1329				
1330				
1331				
1332				
1333				
1334				
1335				
1336				
1337				
1338				
1339				
1340				
1341				
1342	Figure 12: Illustration of FetalMind error samples identified during evaluation.			
1343				
1344				
1345				
1346				
1347				
1348				
1349				

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

Obstetric Ultrasound Report

This report is based on ISUOG and AIUM examination guidelines




I. Patient Information

Name: _____ Record No.: _____ Date of Birth: _____
 Examination Date: _____ Referring Physician: _____ GA by LMP: _____

II. Examination Information

Indication: _____
 Image Quality: _____
 Pregnancy Status: _____

III. Fetal Biometry

Parameter	Measurement (mm)	GA Equivalent (weeks)	Centile
Crown-Rump Length (CRL)	_____	_____	_____
Biparietal Diameter (BPD)	_____	_____	_____
Head Circumference (HC)	_____	_____	_____
Abdominal Circumference (AC)	_____	_____	_____
Femur Length (FL)	_____	_____	_____

Composite GA (US): _____
 Estimated Fetal Weight (Hadlock): _____

IV. Structural Survey

Head: _____ Abdomen: _____
 Face: _____ Spine: _____
 Neck: _____ Limbs: _____
 Thorax: _____ Others: _____
 Heart: _____

V. Conclusion

Complete study, no abnormality detected Incomplete study, no abnormality detected Abnormal findings present

VI. Recommendations

No follow-up required Re-examination after _____ weeks Referral to _____ Others

Sonographer: _____ Signature: _____ Reviewer: _____ Date: _____
 Signature: _____

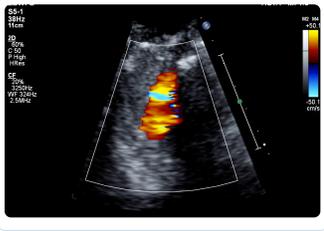
This report was generated by _____ Department of Ultrasound | Image Archive ID: _____

Figure 13: The generalized version of our obstetric ultrasound report template, established with reference to multiple international clinical guidelines. It provides a consistent and clinically grounded format for training and evaluating deep learning systems.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

产科超声检查报告

本报告依据 ISUOG 及 AIUM 检查指南制定




一、患者信息

姓名: _____ 病历号: _____ 出生日期: _____
 检查日期: _____ 转诊医生: _____ 末次月经孕周: _____

二、检查信息

临床指征: _____
 图像质量: _____
 妊娠情况: _____

三、胎儿测量

项目	测量值 (mm)	孕周等效 (周)	百分位 (Centile)
头臀长 (CRL)	_____	_____	_____
双顶径 (BPD)	_____	_____	_____
头围 (HC)	_____	_____	_____
腹围 (AC)	_____	_____	_____
股骨长 (FL)	_____	_____	_____

综合超声孕周: _____
 估计胎重 (Hadlock 法): _____

四、超声可见

头部: _____ 腹部: _____
 面部: _____ 脊柱: _____
 颈部: _____ 四肢: _____
 胸腔: _____ 其他部位: _____
 心脏: _____

五、超声提示

检查完整, 未见异常 检查不全, 未见异常 存在异常发现

六、管理建议

无需复查 周后复查 转诊至 _____ 其他

检查技师: _____ 审核医生: _____
 签名: _____ 日期: _____
 签名: _____

本报告由 _____ 生成 - 超声科 | 图像存储编号: _____

Figure 14: The Chinese version of our obstetric ultrasound report template, established with reference to multiple international clinical guidelines. It provides a consistent and clinically grounded format for training and evaluating deep learning systems.