# MedMO: Grounding and Understanding Multimodal Large Language Model for Medical Images

Ankan Deria*, Komal Kumar*, Adinath Madhavrao Dukre, Eran Segal, Salman Khan, Imran Razzak

Mohamed bin Zayed University of Artificial Intelligence

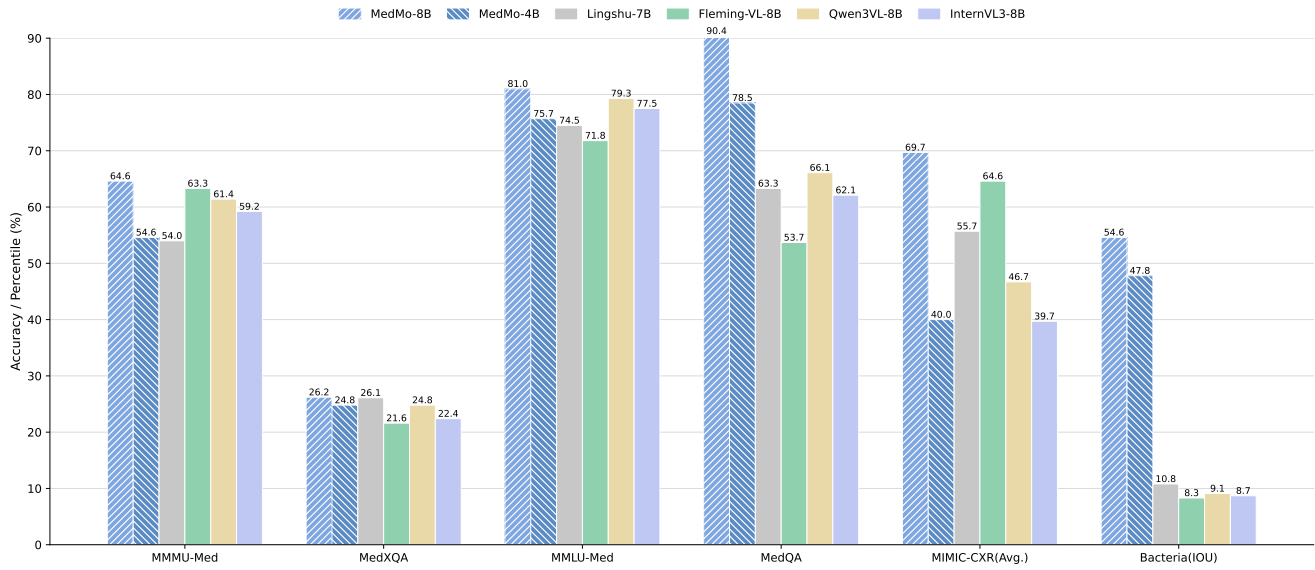{ankan.deria, komal.kumar}@mbzuai.ac.ae

Figure 1. Benchmark performance of **MedMO-8B** across diverse medical imaging tasks, including VQA, QA, report generation, and grounding. MedMO achieves consistent gains over prior models, with improvements of +1.3% on MMMU-Med, +0.1% on MedXQA, +0.7% on MMLU-Med, +24.3% on MedQA, +5.1% on MIMIC-CXR, and a substantial +43.8 IoU on Bacteria segmentation. The large boost in Bacteria IoU stems from the incorporation of fine-grained grounding supervision and high-resolution microscopy data, highlighting MedMO's enhanced spatial reasoning and localization capabilities.

## Abstract

*Multimodal large language models (MLLMs) have rapidly advanced, yet their adoption in medicine remains limited by gaps in domain coverage, modality alignment, and grounded reasoning. In this work, we introduce **MedMO**, a medical foundation model built upon a generalized MLLM architecture and trained exclusively on large-scale, domain-specific data. MedMO follows a multi-stage training recipe: (i) cross-modal pretraining to align heterogeneous visual encoders with a medical language backbone; (ii) instruction tuning on multi-task supervision that spans captioning, VQA, report generation, retrieval, and grounded disease local-ization with bounding boxes; and (iii) reinforcement learning with verifiable rewards that combine factuality checks with a box-level GIoU reward to strengthen spatial ground-ing and step-by-step reasoning in complex clinical scenar-ios. MedMO consistently outperforms strong open-source medical MLLMs across multiple modalities and tasks. On VQA benchmarks, MedMO achieves an average accuracy improvement of **+13.7%** over the baseline and performs within **1.9%** of the SOTA Fleming-VL. For text-based QA, it attains **+6.9%** over the baseline and **+14.5%** over Fleming-VL. In medical report generation, MedMO delivers signif-icant gains in both semantic and clinical accuracy. More-over, it exhibits strong grounding capability, achieving an*

1

*IoU improvement of **+40.4** over the baseline and **+37.0%** over Fleming-VL, underscoring its robust spatial reasoning and localization performance. Evaluations across radiology, ophthalmology, and pathology–microscopy confirm MedMO's broad cross-modality generalization. We release two versions of MedMO: 4B and 8B. Project is available at [genmilab.github.io/MedMO-Page](genmilab.github.io/MedMO-Page)*

## 1. Introduction

Recent advancements in Multimodal Large Language Models (MLLMs) have significantly accelerated progress across multimodal reasoning and understanding tasks [16, 21, 39, 62, 112]. These models unify vision and language comprehension, achieving near-human performance on tasks such as image captioning, visual question answering (VQA), and multimodal reasoning. Despite their remarkable capabilities in general domains, their application to the medical domain remains substantially limited [43, 61, 99]. Biomedical data fundamentally differ from web-scale vision–language pairs: medical images demand precise, domain-specific interpretation, often requiring expert contextualization and robust grounding to textual clinical knowledge [46]. As a result, general-purpose models frequently produce uncertain or hallucinated outputs when applied to medical tasks [14, 45].

To overcome these challenges, recent efforts have sought to adapt general-domain MLLMs into specialized medical multimodal models by incorporating domain-specific data and supervision [32, 76, 93, 107, 110]. Early models such as LLaVA-Med[45] leveraged PubMed-derived datasets for aligning medical images with textual knowledge, achieving foundational progress but limited by noisy data and narrow modality coverage. Subsequent works such as HuatuoGPT-Vision [14], GMAI-VL [46], and HealthGPT [47] introduced high-quality datasets, refined post-training strategies, and scaling recipes that improved alignment and reasoning. Parallel advancements in reasoning-based models, such as OpenAI's o-series [63, 65] and DeepSeek-R1 [23], as well as reinforcement learning with verifiable rewards (RLVR)[77, 104], have inspired recent medical research efforts[40, 68] toward enhancing reasoning reliability and factual grounding in clinical scenarios.

Nevertheless, three key limitations persist in existing medical MLLMs. **(1)** The majority rely on distilled data from advanced proprietary models [21, 22, 62–65], which, while scalable, often lack accurate domain grounding, particularly for fine-grained clinical reasoning. **(2)** Distillation pipelines frequently depend solely on generative outputs without structured supervision, amplifying hallucination risks and inconsistencies. **(3)** Current models focus on individual tasks or narrow modality subsets (e.g., radiology or pathology) rather than achieving unified, cross-modal generalization across the diverse imaging modalities prevalent in real-world healthcare.

To bridge these gaps, we introduce MedMO, a powerful open-source post-trained multimodal large vision–language model (VLM) purpose-built for comprehensive medical image understanding and grounding (See Figure 1). MedMO is developed through a scalable and modular post-training pipeline, emphasizing progressive multimodal alignment, domain-specific reasoning, and cross-modal robustness. We curate and harmonize a 26M+ with 45 open-source multimodal dataset, combining diverse medical imaging modalities (radiology, pathology, ophthalmology, dermatology, CT, MRI, ultrasound, and surgical videos) with carefully aligned text sources from open biomedical corpora and general-domain visual data. Through multi-stage post-training, MedMO progressively enhances its capacity for visual grounding, clinical reasoning, and textual alignment, establishing a scalable pipeline toward a generalist foundation multimodal model for medical AI.

We further conduct comprehensive experiments and analyses on data curation, training, and alignment strategies, providing a transparent and reproducible framework for future medical MLLM development. Extensive evaluations demonstrate that MedMO achieves state-of-the-art (SOTA) performance across diverse benchmarks, surpassing prior open and proprietary systems on tasks including medical VQA, report generation, and diagnostic reasoning.

Our main **contributions** are summarized as follows:
- We develop a powerful open-source post-trained multimodal large VLM, **MedMO**, designed for comprehensive medical image understanding and grounding.
- We curate over 26M multimodal medical and biomedical samples from **45 datasets** and establish a multi-stage post-training that progressively enhances cross-modal alignment and reasoning. This provides a scalable roadmap toward a generalist foundation model for medical.
- To evaluate VLM performance on detection tasks, we construct a dedicated Cell dataset from opensource microscopy images with varying sizes, shapes, and densities.
- We conduct extensive experiments and analyses across data and methodology dimensions, providing an open benchmark for future multimodal medical LLM research and training recipes.

## 2. Related Works

### 2.1. Medical Language Multi-model Models

The rapid progress of LLMs has catalyzed remarkable advances in medical images. Building upon the success of general-domain LLMs, researchers have developed domain-adapted medical MLLMs that integrate visual and textual reasoning for clinical understanding [4, 84]. Early efforts such as LLaVA-Med [45], Med-Flamingo [58], Qilin-MedVL [50], and BioMedGPT [107] established the first

medical vision–language models by aligning specialized visual encoders with pre-trained LLMs via linear projection layers, enabling foundational multimodal reasoning. However, these early systems were constrained by limited data diversity and suboptimal modality alignment, leading to hallucinations and factual inconsistencies [14, 45]. Subsequent studies expanded this paradigm through richer datasets [30, 33, 46], improved training strategies [60, 91], efficient fine-tuning [47], and reinforcement learning [40, 68]. Proprietary systems such as Med-Gemini[20] and Med-PaLM [80, 81] have further integrated multimodal and structured data for advanced reasoning, achieving strong performance across diagnostic and question-answering tasks [5, 6, 73, 95, 101]. Concurrently, specialized MLLMs targeting specific clinical contexts such as pathology [53, 76, 92, 111], radiology [19, 32, 66, 79, 83, 106], and ophthalmology [25] have emerged, highlighting the growing demand for fine-grained, modality-aware intelligence in medical. Recent SOTA frameworks, such as Lingshu [97] and Fleming-vl [78], have improved the integration of medical and natural VLM tasks. However, their capabilities remain limited to selective tasks. Building on these foundations, our work emphasizes large-scale open-source post-training and progressive multimodal alignment. MedMO adopts a multi-stage design leveraging over 26M diverse multimodal samples, unifying heterogeneous medical modalities and textual data to achieve substantial gains across diverse clinical tasks.

## 2.2. Grounding using multi-model models

Unlike detection objective-based approaches such as grounding-DINO [51], recent flagship VLMs have moved beyond captioning/VQA to explicit visual grounding [18, 21, 51] as well as point grounding [24], i.e., returning spatial evidence such as bounding boxes or points aligned to textual queries. The Qwen2.5-VL [8] report highlights grounding as a built-in capability, emphasizing precise object localization and event localization in long videos through native dynamic-resolution processing and absolute time encoding. Qwen2.5-VL generates grounded outputs in JSON with absolute coordinates, supporting both boxes and point clicks [24]. Although the technical report is general-domain, these grounding primitives transfer to clinical data. For instance, MedSG-Bench [105] evaluates sequential medical grounding (difference/consistency grounding across image series) and explicitly benchmarks Qwen2.5-VL alongside medical-domain MLLMs (e.g., HuatuoGPT-Vision [15]), finding that even advanced VLMs still face challenges on fine-grained, clinically realistic localization tasks-underscoring the need for domain-aligned post-training.

## 3. Methodology-MedMO

The overall methodology and multi-stage training pipeline are provided in Figure 2. Starting from the Qwen3-VL-8B-Instruct model[1], our approach consists of four sequential post-training stages: (1) General SFT aimed to train on large-scale instruction data to build foundational medical understanding; (2) High-quality medical image supervised fine-tuning, focused on expert-curated data to enhance visual grounding; (3) Instruction tuning and grounding fine-tuning, which align the model with clinical answering and spatial localization tasks; and (4) Reinforcement learning, designed to further improve instruction-following behavior and grounding accuracy. The following subsections provide an overview of the supervised fine-tuning strategy, followed by detailed descriptions of each stage.

## 3.1. Overview of Supervised Fine-tuning

Our supervised fine-tuning (SFT) approach follows the standard next-token prediction paradigm for vision-language models. Given a multimodal input consisting of an image $\mathbf{v}$ and text sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, the model learns to predict the target response $\mathbf{y} = \{y_1, y_2, \ldots, y_m\}$ by maximizing the conditional likelihood:

$$\mathcal{L}_{\text{SFT}} = -\sum_{i=1}^{m} \log p_\theta(y_i \mid \mathbf{v}, \mathbf{x}, y_{<i}), \quad (1)$$

where $\theta$ represents the model parameters, and $y_{<i}$ denotes all previously generated tokens. MedMO builds upon the Qwen3-VL architecture, which consists of three primary components: (1) a vision encoder $\mathcal{E}_v$ that processes input images into visual representations; (2) a vision–language adapter $\mathcal{A}$ that projects multi-level ViT features into the language model's embedding space through a DeepStack fusion mechanism, capturing fine-grained visual details and enhancing image–text alignment; and (3) a large language model decoder $\mathcal{D}$ that generates textual responses.

## 3.2. Stage 1: General Medical SFT

The first stage aims to establish foundational medical knowledge across diverse modalities and clinical scenarios. We utilize the publicly available **MedTrinity** dataset [96], comprising **18.5M** large-scale instruction-following samples. This dataset $\mathcal{D}_{\text{general}}$ spans multiple imaging modalities (X-ray, CT, MRI, ultrasound, pathology, etc.) and includes captioning, visual question answering (VQA), and general-domain multimodal tasks, as illustrated in Figure 4.

The Stage 1 dataset consists of:

- **Medical image captioning:** $\mathcal{D}_{\text{caption}}$ with detailed textual descriptions of medical images.
- **Medical VQA:** $\mathcal{D}_{\text{vqa}}$ covering disease identification, anatomical recognition, and reasoning tasks.
- **General multimodal data:** $\mathcal{D}_{\text{general-mm}}$ for maintaining broad visual–language alignment.
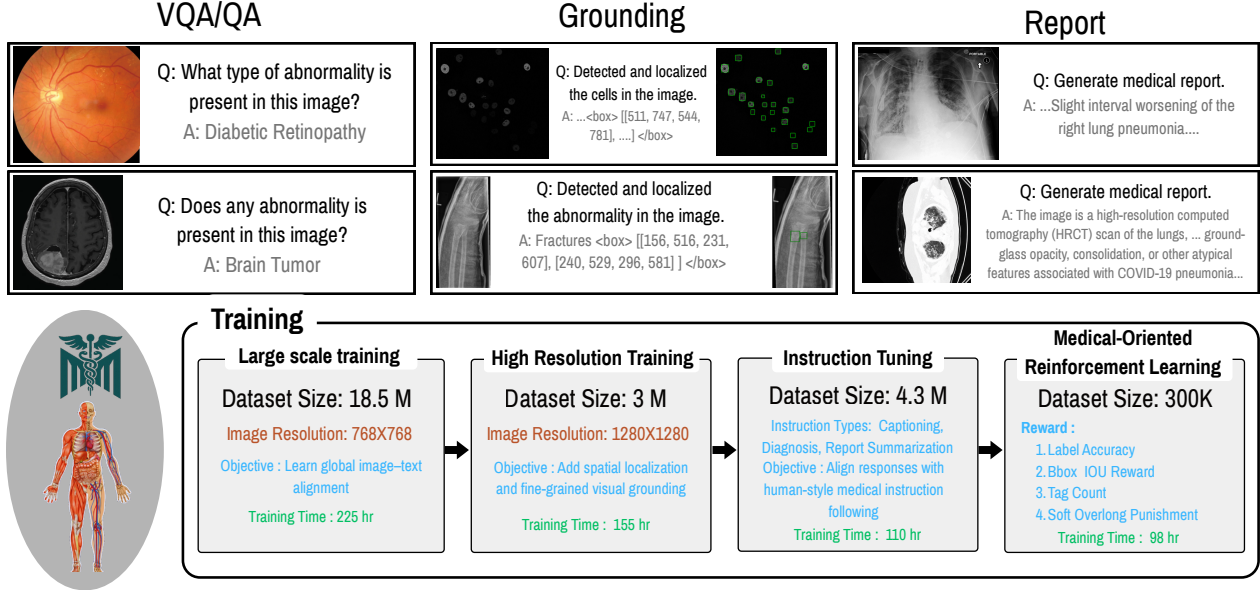
---

[1]Qwen/Qwen3-VL-8B-Instruct

3

Figure 2. **Overview of the multi-stage training pipeline for medical image analysis.** The workflow consists of three main capabilities: (Top row) VQA/QA for identifying abnormalities in medical images, Grounding for spatial localization of detected features with bounding box coordinates, and Report generation for producing detailed medical reports. (Bottom) The training pipeline progresses through four sequential stages: (1) *Large-scale training* on 18.5M image-text pairs at 768×768 resolution for global image-text alignment, (2) *High-resolution training* on 3M samples at 1280×1280 resolution to enhance spatial localization and fine-grained visual grounding, (3) *Instruction tuning* on 4.3M samples covering captioning, diagnosis, and report summarization tasks to align responses with human-style medical instruction following, and (4) *Medical-oriented reinforcement learning* on 300K samples optimized using four reward signals: label accuracy, bounding box IoU, tag count, and soft overlap punishment. The complete pipeline for the MedMO-8B.

The combined dataset is defined as:

$$\mathcal{D}_{\text{stage1}} = \mathcal{D}_{\text{caption}} \cup \mathcal{D}_{\text{vqa}} \cup \mathcal{D}_{\text{general-mm}}. \tag{2}$$

### 3.3. Stage 2: Quality Medical Image and Grounding

The second stage of SFT focuses on high-quality, expert-annotated medical image–text pairs to strengthen visual understanding and introduce grounding capability. We curate a refined dataset $\mathcal{D}_{\text{hq}}$ that includes both standard image–text supervision and medical grounding datasets containing bounding-box annotations (e.g., Chest X-ray, Wrist X-ray, Cell Microscopy, and CT). This stage extends the model's visual encoder to predict localized features and bounding box coordinates, enabling spatial awareness while preserving global image–text alignment. Training objectives remain consistent with Stage 3.2, combining captioning and VQA with supervised grounding signals.

**Grounding Dataset.** The grounding dataset $\mathcal{D}_{\text{ground}}$ includes: (1) Object detection annotations for anatomical structures and lesions, (2) Referring expression comprehension, and (3) Visual grounding QA pairs for spatial localization.

### 3.4. Stage 3: Instruction Tuning

The third stage aligns MedMO's responses with human-style medical reasoning through instruction tuning. Using a dataset $\mathcal{D}_{\text{inst}}$ of 4.3M multimodal instruction–response pairs, this phase covers captioning, diagnostic question answering, report summarization, and retrieval-based reasoning tasks. Instruction tuning improves task generalization and factual consistency, integrating clinical context understanding into both text- and vision-guided reasoning.

### 3.5. Stage 4: Reinforcement Learning

The final stage employs GRPO [77] to enhance instruction-following capabilities through preference learning.

**GRPO Objective.** It optimizes the model by comparing multiple sampled responses for the same input. For each input $(\mathbf{v}, \mathbf{x})$, we sample $G$ responses $\{\mathbf{y}^{(1)}, \ldots, \mathbf{y}^{(G)}\}$ from the current policy $\pi_\theta$. Each response is evaluated using a reward function $r(\mathbf{v}, \mathbf{x}, \mathbf{y})$ that measures quality.

We follow the same objective as in GRPO [23, 77] with clip-higher and token level loss motivated from from DAPO [104]. For $(q, a) \sim \mathcal{D}$, $\{o_i\}_{i=1}^{G} \sim \pi_{\theta_{\text{old}}}(\cdot|q)$,

$$
\begin{aligned}
J(\theta) =& \mathbb{E}_{(q,a),o_i} \left[ \frac{1}{\sum_{i=1}^{G} |o_i|} \sum_{i=1}^{G} \sum_{t=1}^{|o_i|} \min \left( r_{i,t}(\theta) \hat{A}_{i,t}, \right. \right. \\
& \left. \left. \text{clip} \left( r_{i,t}(\theta), 1 - \varepsilon_{\text{low}}, 1 + \varepsilon_{\text{high}} \right) \hat{A}_{i,t} \right) \right]
\end{aligned} \tag{3}
$$

4

$$r_{i,t}(\theta) = \frac{\pi_\theta(o_{i,t} \mid q, o_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t} \mid q, o_{i,<t})}, \tag{4}$$

$$\hat{A}_{i,t} = \frac{R_i - \text{mean}(\{R_i\}_{i=1}^G)}{\text{std}(\{R_i\}_{i=1}^G)}. \tag{5}$$

The KL divergence term ensures the policy doesn't deviate too far from the reference model $\pi_{\text{ref}}$:

$$\mathcal{L}_{\text{KL}} = \mathbb{E}_{(\mathbf{v},\mathbf{x},\mathbf{y})} \left[ D_{\text{KL}}(\pi_\theta(\cdot \mid \mathbf{v}, \mathbf{x}) \| \pi_{\text{ref}}(\cdot \mid \mathbf{v}, \mathbf{x})) \right]. \tag{6}$$

For the reward function, we combine label accuracy, bounding-box reward, tag count, and soft-overlap penalty (see Fig. 2). While these components are common in RL-based training, we introduce the Bounding Box Reward as a verifiable, spatially grounded signal that directly enhances localization performance.

### 3.5.1. Bounding Box Reward

Given ground truth boxes $\mathcal{G} = \{g_j\}_{j=1}^G$ and predictions $\mathcal{P} = \{p_i\}_{i=1}^P$ in XYXY format and $\text{GIoU}_{ij} \in [-1, 1]$ [71], we score pairs via

$$L1_{ij} = \frac{|x_1^p - x_1^g| + |y_1^p - y_1^g| + |x_2^p - x_2^g| + |y_2^p - y_2^g|}{2\sqrt{H^2 + W^2}}$$

Normalize by the average image dimension makes the denominator resolution-invariant and proportional to image diagonal length. We obtain a one-to-one assignment $M \subseteq \{1 \ldots P\} \times \{1 \ldots G\}$ by Hungarian matching on

$$C_{ij} = w_{\text{L1}}^m L1_{ij} + w_{\text{G}}^m (1 - \text{GIoU}_{ij}), \quad w_{\text{L1}}^m = 5, \ w_{\text{G}}^m = 2.$$

For each matched pair $(i, j) \in M$, define a per-pair quality

$$s_{ij} = \frac{w_{\text{L1}} (1 - \text{clip}_{[0,1]}(L1_{ij})) + w_{\text{G}} \left(\frac{\text{GIoU}_{ij}+1}{2}\right)}{w_{\text{L1}} + w_{\text{G}}},$$

where $w_{\text{L1}} = 5$, $w_{\text{G}} = 2$. The reward is a coverage-normalized sum with optional FP/FN penalties (Pen):

$$B = \frac{1}{G} \sum_{(i,j) \in M} s_{ij}, \ \text{Pen} = \frac{\lambda_{\text{FN}}(G - |M|) + \lambda_{\text{FP}}(P - |M|)}{\max(1, G)},$$

$$\boxed{R_{\text{bbox}} = \text{clip}_{[0,1]}(B - \text{Pen})}^2.$$

## 4. Experiments

### 4.1. Experimental Setup

MedMO was trained using **64× AMD Instinct MI210 GPUs** (64 GB each) for **25 days** following a four-stage progressive pipeline (Figure 2). The stages comprised: large-scale general medical SFT on 18.5M image–text pairs at 768×768 resolution (**225 h**); high-resolution fine-tuning on

3M curated samples at 1280×1280 (**155 h**); instruction tuning on 4.3M multimodal examples covering captioning, diagnosis, and report summarization (**110 h**); and medical-oriented reinforcement learning on 300K samples with rewards for label accuracy, bounding-box IoU (**98 h**). We follow standard VLM training practices using TRL [88]. Stage 1 uses BS = 10, LR = 1e-5, cosine schedule, and grad accum = 2. Stage 2 adopts BS = 2, LR = 8e-6, cosine schedule, and grad accum = 8. Stage 3 employs BS = 10 and LR = 5e-6 with grad accum = 2 for stable convergence [2].

### 4.2. Datasets

We assembled a unified multimodal corpus of **45 datasets** spanning radiology, pathology, ophthalmology, dermatology, and surgical imaging, totaling over **26M samples**. The **MedTrinity** dataset [96] forms the core, contributing **18.5M** public instruction-following pairs. The corpus combines image–text and text-only data across diverse medical domains and clinical tasks. The dataset (Figure 4) covers both *imaging modalities* (e.g., X-ray, CT, MRI, ultrasound, optical, and nuclear imaging) and *biological systems* (chest, brain, heart, liver, kidney, eye, colon, and tissue). For grounding tasks, we additionally used datasets with bounding-box annotations, including *Chest X-ray, Wrist X-ray, Cell microscopy, and CT* images. This comprehensive coverage supports robust multimodal understanding, spatial reasoning, and medical grounding. We curate a Cell Benchmark Dataset from open-source microscopy images[2], such as DeepCell [9] and Bacteria [87], covering diverse cell counts and densities[2].

### 4.3. Results and Analysis

#### 4.3.1. SOTA comparison of MedMO for QnA

Table 1 summarizes MedMO's performance across medical VQA and Text QA benchmarks on two variants $4B$ & $8B$. MedMO achieves state-of-the-art results, matching specialized medical models such as Fleming-VL-8B (64.4%) and Lingshu-7B (61.8%), and reaching the highest score on MMMU-Med (64.6%).

On Text QA benchmarks, MedMO achieves an average accuracy of **61.4%**, surpassing Qwen3VL-8B (54.5%) by 7%. The improvement is most notable on reasoning tasks such as MedQA (90.4%), MMLU-Med (81.0%), and MedMCQA (65.0%), demonstrating MedMO's enhanced medical knowledge integration and clinical reasoning capability.

#### 4.3.2. SOTA comparison of MedMO for understanding

Table 2 evaluates medical report generation capabilities across four datasets using both semantic-based metrics (ROUGE-L, CIDEr) and model-based metrics (RaTE, Semb).

**MIMIC-CXR Results.** On MIMIC-CXR, the most widely-used benchmark for chest X-ray report generation, our model

---

[2]For more details, please see our Supplementary Material.

Table 1. Performance comparison across medical **VQA** and **Text QA** benchmarks. **Bold** and underline indicate the best and second-best results, respectively. OMIVQA and MedXQA refer to the OmniMedVQA and MedXpertQA benchmarks.

| Models | VQA Benchmarks | | | | | | | | Text QA Benchmarks | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | MMMU-Med | VQA-RAD | SLAKE | PathVQA | PMC-VQA | OMVQA | MedXQA | Avg. | MMLU-Med | PubMedQA | MedMCQA | MedQA | Medbullets | MedXQA | SGPQA | Avg. |
| *Closed-source Models* | | | | | | | | | | | | | | | | |
| GPT-4.1 | 75.2 | 65.0 | 72.2 | 55.5 | 55.2 | 75.5 | 45.2 | 63.4 | 89.6 | 75.6 | 77.7 | 89.1 | 77.0 | 30.9 | 49.9 | 70.0 |
| Claude Sonnet 4 | 74.6 | 67.6 | 70.6 | 54.2 | 54.4 | 65.5 | 43.3 | 61.5 | 91.3 | 78.6 | 79.3 | 92.1 | 80.2 | 33.6 | 56.3 | 73.1 |
| Gemini-2.5-Flash | 76.9 | 68.5 | 75.8 | 55.4 | 55.4 | 71.0 | 52.8 | 65.1 | 84.2 | 73.8 | 73.6 | 91.2 | 77.6 | 35.6 | 53.3 | 69.9 |
| *Open-source Models* | | | | | | | | | | | | | | | | |
| BiomedGPT | 24.9 | 16.6 | 13.6 | 11.3 | 27.6 | 27.9 | – | – | – | – | – | – | – | – | – | – |
| Med-R1-2B | 34.8 | 39.0 | 54.5 | 15.3 | 47.4 | – | 21.1 | – | 51.5 | 66.2 | 39.1 | 39.9 | 33.6 | 11.2 | 17.9 | 37.0 |
| MedVLM-R1-2B | 35.2 | 48.6 | 56.0 | 32.5 | 47.6 | 77.7 | 20.4 | 45.4 | 51.8 | 66.4 | 39.7 | 42.3 | 33.8 | 11.8 | 19.1 | 37.8 |
| MedGemma-4B-IT | 43.7 | 72.5 | 76.4 | 48.8 | 49.9 | 69.8 | 22.3 | 54.8 | 66.7 | 72.2 | 52.2 | 56.2 | 45.6 | 12.8 | 21.6 | 46.8 |
| LLaVA-Med-7B | 29.3 | 53.7 | 48.0 | 38.8 | 30.5 | 44.3 | 20.3 | 37.8 | 50.6 | 26.4 | 39.4 | 42.0 | 34.4 | 9.9 | 16.1 | 31.3 |
| HuatuoGPT-V-7B | 47.3 | <u>67.0</u> | 67.8 | 48.0 | 53.3 | 74.2 | 21.6 | 54.2 | 69.3 | 72.8 | 51.2 | 52.9 | 40.9 | 10.1 | 21.9 | 45.6 |
| BioMediX2-8B | 39.8 | 49.2 | 57.7 | 37.0 | 43.5 | 63.3 | 21.8 | 44.6 | 68.6 | 75.2 | 52.9 | 58.9 | 45.9 | 13.4 | 25.2 | 48.6 |
| Qwen2.5VL-7B | 50.6 | 64.5 | 67.2 | 44.1 | 51.9 | 63.6 | 22.3 | 52.0 | 73.4 | 76.4 | 52.6 | 57.3 | 42.1 | 12.8 | 26.3 | 48.7 |
| InternVL2.5-8B | 53.5 | 59.4 | 69.0 | 42.1 | 51.3 | 81.3 | 21.7 | 54.0 | 74.2 | 76.4 | 52.4 | 53.7 | 42.4 | 11.6 | 26.1 | 48.1 |
| InternVL3-8B | 59.2 | 65.4 | 72.8 | 48.6 | 53.8 | 79.1 | 22.4 | 57.3 | 77.5 | 75.4 | 57.7 | 62.1 | 48.5 | 13.1 | 31.2 | 52.2 |
| Lingshu-7B | 54.0 | 67.9 | <u>83.1</u> | <u>61.9</u> | 56.3 | 82.9 | <u>26.1</u> | 61.8 | 74.5 | 76.6 | 55.9 | 63.3 | 56.2 | <u>16.5</u> | 26.3 | 52.8 |
| Fleming-VL-8B | <u>63.3</u> | 66.1 | **86.5** | **62.9** | **64.3** | **86.7** | 21.6 | **64.4** | 71.8 | 74.0 | 51.8 | 53.7 | 40.5 | 12.1 | 24.4 | 46.9 |
| Qwen3VL-8B | 61.4 | 64.1 | 47.3 | 14.6 | 52.3 | 77.2 | 24.8 | 48.8 | <u>79.3</u> | 70.4 | <u>60.0</u> | <u>66.1</u> | 56.1 | 15.1 | <u>34.7</u> | 54.5 |
| MedMO-4B | 54.6 | 50.9 | 41.0 | 62.4 | 50.6 | 79.7 | 24.8 | 52.0 | 75.7 | **78.0** | 58.0 | 78.5 | <u>57.5</u> | 16.4 | 29.4 | <u>56.2</u> |
| MedMO-8B | **64.6** | 64.7 | 81.6 | 56.3 | <u>59.4</u> | <u>84.8</u> | **26.2** | <u>62.5</u> | **81.0** | <u>77.6</u> | **65.0** | **90.4** | **60.2** | **19.9** | **36.0** | **61.4** |

Table 2. Comparison of medical report generation performance on MIMIC-CXR, CheXpert Plus, IU-Xray, and Med-Trinity using semantic (ROUGE-L, CIDEr) and model-based (RaTE, Semb) metrics. Models highlighted in green denote our proposed MedMO, which achieves the best overall performance across all datasets.

| Models | MIMIC-CXR | | | | CheXpert Plus | | | | IU-Xray | | | | Med-Trinity | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ROUGE-L | CIDEr | RaTE | Semb | ROUGE-L | CIDEr | RaTE | Semb | ROUGE-L | CIDEr | RaTE | Semb | ROUGE-L | CIDEr | RaTE | Semb |
| *Closed-source Models* | | | | | | | | | | | | | | | | |
| GPT-4.1 | 9.0 | 82.8 | 51.3 | 23.9 | 24.5 | 78.8 | 45.5 | 23.2 | 30.2 | 124.6 | 51.3 | 47.5 | – | – | – | – |
| Claude Sonnet 4 | 20.0 | 56.6 | 45.6 | 19.7 | 22.0 | 59.5 | 43.5 | 18.9 | 25.4 | 88.3 | 55.4 | 41.0 | – | – | – | – |
| Gemini-2.5-Flash | 25.4 | 80.7 | 50.3 | 29.7 | 23.6 | 72.2 | 44.3 | 27.4 | 33.5 | 129.3 | 55.6 | 50.9 | – | – | – | – |
| *Open-source Models* | | | | | | | | | | | | | | | | |
| Med-R1-2B | 19.3 | 35.4 | 40.6 | 14.8 | 18.6 | 37.1 | 38.5 | 17.8 | 16.1 | 38.3 | 41.4 | 12.5 | – | – | – | – |
| MedVLM-R1-2B | 20.3 | 40.1 | 41.6 | 14.2 | 20.9 | 43.5 | 38.9 | 15.5 | 22.7 | 61.1 | 46.1 | 22.7 | – | – | – | – |
| MedGemma-4B-IT | 25.6 | 81.0 | 52.4 | 29.2 | **27.1** | 79.0 | <u>47.2</u> | 29.3 | 30.8 | 103.6 | 57.0 | 46.8 | – | – | – | – |
| LLaVA-Med-7B | 15.0 | 43.4 | 42.8 | 18.3 | 18.4 | 45.5 | 38.8 | 23.5 | 18.8 | 68.2 | 40.9 | 16.0 | – | – | – | – |
| HuatuoGPT-V-7B | 23.4 | 69.5 | 48.9 | 20.0 | 21.3 | 64.7 | 44.2 | 19.3 | 29.6 | 104.3 | 52.9 | 40.7 | – | – | – | – |
| BioMediX2-8B | 20.0 | 52.8 | 44.4 | 17.7 | 18.1 | 47.9 | 40.8 | 21.6 | 19.6 | 58.8 | 40.1 | 11.6 | – | – | – | – |
| Qwen2.5VL-7B | 24.1 | 63.7 | 47.0 | 18.4 | 22.2 | 62.0 | 41.0 | 17.2 | 26.5 | 78.1 | 48.4 | 36.3 | <u>23.5</u> | 81.5 | 44.9 | <u>38.3</u> |
| InternVL2.5-8B | 23.2 | 61.8 | 47.0 | 21.0 | 20.6 | 58.5 | 43.1 | 19.7 | 24.8 | 75.4 | 51.1 | 36.7 | 13.5 | 47.1 | 42.5 | 12.8 |
| InternVL3-8B | 22.9 | 66.2 | 48.2 | 21.5 | 20.9 | 65.4 | 44.3 | 25.2 | 22.9 | 76.2 | 51.2 | 31.3 | 12.9 | 46.6 | 42.2 | 3.7 |
| Lingshu-7B | 30.8 | 109.4 | 52.1 | 30.0 | <u>26.5</u> | 79.0 | 45.4 | 26.8 | <u>41.2</u> | <u>180.7</u> | <u>57.6</u> | <u>48.4</u> | 16.0 | 74.5 | 44.4 | 24.0 |
| Fleming-VL-8B | **35.7** | <u>132.5</u> | <u>56.7</u> | <u>33.6</u> | 26.1 | <u>82.2</u> | 47.1 | <u>40.1</u> | **44.9** | **198.6** | **66.0** | **51.3** | 13.1 | 35.8 | 41.9 | 18.1 |
| Qwen3VL-8B | 25.1 | 77.9 | 50.3 | 33.4 | 21.9 | 67.4 | 44.4 | 37.9 | 25.0 | 91.44 | 52.5 | 42.9 | 20.2 | 69.9 | 45.9 | 33.6 |
| MedMO-4B | 26.0 | 92.6 | 49.8 | 31.6 | 15.1 | 62.3 | 36.6 | 34.2 | 26.6 | 94.0 | 42.1 | 41.3 | 22.5 | <u>152.6</u> | <u>47.8</u> | 34.3 |
| MedMO-8B | <u>31.7</u> | **140.0** | **57.1** | **50.0** | 23.6 | **87.5** | **47.3** | **42.2** | 31.1 | 169.7 | 45.3 | 41.3 | **37.0** | **270.4** | **53.0** | **39.2** |

achieves outstanding results with a CIDEr score of 140.0 and ROUGE-L of 31.7%, substantially outperforming strong medical baselines including Fleming-VL-8B (132.5, 35.7%) and Lingshu-7B (109.4, 30.8%). While Fleming-VL-8B achieves a higher ROUGE-L, our superior CIDEr score indicates better semantic coherence and clinical relevance. Most notably, the RaTE score of 57.1% and remarkable Semb score of 50.0% demonstrate that MedMO generates reports with high clinical accuracy and semantic fidelity.

**CheXpert Plus and IU-Xray Results.** On CheXpert Plus, MedMO achieves the highest CIDEr (87.5) and strong Semb performance (42.2), though MedGemma-4B-IT achieves slightly higher ROUGE-L (27.1% vs. 23.6%). This suggests different models may excel at different aspects of report generation-lexical overlap versus semantic coherence and clinical accuracy. For IU-Xray, our model achieves strong

CIDEr performance (169.7), second only to Fleming-VL-8B (198.6), demonstrating consistent capability across different chest X-ray datasets. However, the lower RaTE score (45.3%) compared to other datasets suggests this benchmark may have different annotation characteristics or evaluation criteria that warrant further investigation.

**Med-Trinity Results.** Most notably, on the Med-Trinity dataset which encompasses diverse medical imaging modalities beyond radiology, MedMO-8B dramatically outperforms all baselines with a CIDEr score of 270.4-nearly 4× higher than the next best model (Qwen2.5VL-7B at 81.5). This substantial margin, combined with the highest ROUGE-L (37.0%) and RaTE (53.0%), demonstrates exceptional capability in handling diverse medical imaging modalities, including CT, MRI, ultrasound, and pathology images.

**Question: What is the name of the skin abnormality in this image?**
Options : A. Eczema, B. Squamous cell carcinoma, C. Malignant melanoma, D. Melanoma

**Ground Truth:** B. Squamous cell carcinoma
**Lingshu:** B. Psoriasis
**Fleming VL-8B:** B. Psoriasis
**Qwen3-VL :** B. Psoriasis
**MedMO:** B. Squamous cell carcinoma

**Question: Detect and localize all cell.**
**Ground Truth :** [[559,447,587,477], [583,437,603,471]]

**Lingshu :** [448,502,500,542], [448,542,500,435584]
**Fleming VL :** [496,391,566,446]
**Qwen3 VL :** [588,396,642,440]
**MedMO:** Bacteria: [[548,441,594,470], [578,440,609,465]]

**Question : Detect and localize all cell.**
**Ground Truth :** [54,545,63,554]

**Lingshu :** [0,0,644,644]
**Fleming VL :** [0,0,999,999]
**Qwen3 VL :** [31,21,965,957]
**MedMO:** Cell: [53,548,62,557]

**Groundtruth:** Findings: There is still a small right apical pneumothorax, unchanged from 11:26 a.m. There is no shift of the mediastinum. The lungs are clear. The cardiomediastinal silhouette is within normal limits. Impression: Stable right apical pneumothorax.

**Lingshu:** Findings: AP upright and lateral views of the chest were provided. The lungs are clear. No large effusion or pneumothorax. The heart size is normal. Mediastinal contour is unremarkable. Bony structures are intact. No free air below the right hemidiaphragm. Clips are noted in the right upper quadrant. Degenerative changes are noted in the spine. Impression: No acute findings.

**Fleming-VL :** Findings: The lungs are clear. The heart size is normal. The mediastinal contours are normal. There are no pleural effusions. No pneumothorax is seen. The bony thorax is grossly intact. Impression: No acute cardiac or pulmonary process.

**Qwen3-VL :** The chest X-ray is an anteroposterior (AP) view, as indicated by the "AP" label and the "L" marker pointing to the left side. The bony structures, including the ribs, clavicles, and spine, appear intact without obvious fractures or dislocations. The lung fields are generally clear, with no obvious consolidation, pleural effusion, or pneumothorax. The cardiac silhouette is within normal size limits, and the diaphragm appears well-defined. There are no visible foreign bodies or significant abnormalities in the mediastinum. The soft tissues of the chest wall appear normal.

**MedMO:** The image is a chest X-ray showing the thoracic cavity with the lungs, heart, and the bony structures of the chest wall. The lungs are the primary organs visible, with the heart situated medially between them. The region of interest, located in the upper right portion of the image, shows a translucent area without lung markings, indicative of pneumothorax. This area is adjacent to the pleural line, which is abnormal as it suggests the presence of air in the pleural space, causing lung collapse. The relative position of this region to the rest of the lung tissue is critical, as it can affect the surrounding lung areas and potentially lead to a shift in the mediastinum or trachea if it progresses to a tension pneumothorax.
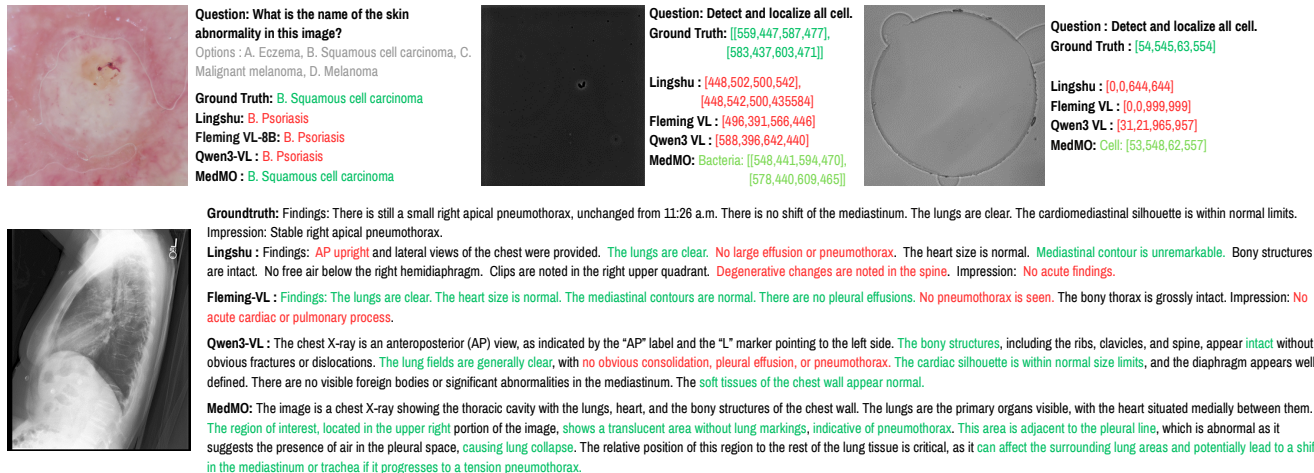
Figure 3. Qualitative comparison across diverse medical and visual question-answering tasks. Each block shows the ground truth, model predictions from Fleming-VL-8B (current Medical SOTA), Qwen3-VL (Baseline), and MedMO, and highlights textual or spatial alignment. MedMO provides more accurate medical understanding and localization in both diagnostic accuracy and clinical reasoning.

Figure 4. Composition of the unified multi-modal medical dataset comprising diverse imaging modalities and biological systems.

Table 3. Performance comparison of selected MLLMs on Medical Grounding Benchmarks. NIH: Chest X-ray; DeepLesion: lesion detection; Bacteria: detection; MedSG: multi-view, object-tracking, and referring tasks. All values are IoU scores (%), and "Avg." denotes the mean across tasks.

| Model | NIH | DeepLession | Bacteria | MedSG (multi_view) | MedSG (object_tracking) | MedSG (referring) | Avg. |
|---|---|---|---|---|---|---|---|
| InternVL3-8B | 10.1 | 0.00 | 0.7 | 6.3 | 13.0 | 3.3 | 5.6 |
| Fleming-VL-8B | 0.00 | 0.00 | 8.3 | 42.0 | 36.7 | 16.6 | 17.2 |
| Lingshu-7B | 5.3 | 0.7 | 0.00 | 28.3 | 38.7 | 10.4 | 13.9 |
| Qwen3VL-8B | **16.4** | 0.00 | 9.16 | 8.4 | 17.8 | 31.4 | 13.8 |
| MedSG-Bench | – | – | – | 55.0 | 62.1 | 60.4 | - |
| MedMO-8B | 8.83 | **38.5** | **54.6** | **75.8** | **77.2** | **70.1** | **54.2** |

ment suggests our training methodology effectively addresses the grounding capabilities often lacking in general-purpose vision-language models. For bacteria segmentation, a challenging fine-grained task, MedMO achieves 54.6% IoU, significantly outperforming Qwen3VL-8B (9.16%) and Fleming-VL-8B (8.3%). On NIH Chest X-ray, while Qwen3VL-8B achieves the highest score (16.4%), our model demonstrates competitive performance (8.83%), indicating reasonable localization ability on this challenging dataset with pathological findings.

**MedSG Multi-task Benchmarks.** On the multi-task MedSG benchmarks, MedMO-8B shows remarkable performance gains: 75.8% on multi-view tasks (requiring correspondence across different imaging views), 77.2% on object tracking (temporal consistency), and 70.1% on referring expression grounding (natural language-based localization). These results substantially exceed the specialized MedSG-Bench model (55.0%, 62.1%, 60.0%), demonstrating that our general-purpose medical model successfully learns fine-grained spatial understanding without sacrificing performance on other tasks. The consistent improvements

### 4.3.3. MedMO for grounding

Table 3 presents performance on medical grounding benchmarks measured by Intersection over Union (IoU). MedMO-8B achieves state-of-the-art results across all tasks, demonstrating strong spatial reasoning and localization capabilities crucial for medical image analysis.

**Lesion Detection.** On DeepLesion, our model achieves 38.5% IoU, substantially outperforming all baselines that either fail completely (0.00%) or show minimal performance (Lingshu-7B at 0.7%). This significant improve-

across diverse grounding tasks-from anatomical structure localization to microscopic bacteria segmentation to multi-view reasoning-highlight the robustness and generalizability of MedMO's visual encoding and grounding mechanisms.

### 4.3.4. MedMO improves medical analysis

MedMO demonstrates superior reasoning and localization across diverse medical imaging tasks, including dermatology, radiology, and cellular detection. As shown in Fig. 3, MedMO not only aligns more closely with clinical ground truth but also provides coherent textual explanations and accurate spatial grounding. Compared to Fleming-VL-8B and Qwen3-VL, MedMO exhibits stronger diagnostic understanding, improved interpretability, and consistent bounding-box precision, reflecting its capability to bridge medical vision and language reasoning effectively.

## 4.4. Ablation Study

### 4.4.1. Ablation on Post-Training Stages

We perform a stage-wise ablation to evaluate the contribution of each post-training phase to MedMO's performance on radiology and QA benchmarks. As shown in Figures 5 and 6, performance progressively improves across stages, validating the effectiveness of our optimization strategy. In Stage 1, the model trained on the MedTrinity dataset achieves strong accuracy on that dataset but shows slight degradation on others. Stage 2 provides the largest gain through high-resolution and diverse medical data training, while Stage 3 further boosts QA and VQA performance via instruction tuning, enhancing multimodal alignment and reasoning. Each stage contributes complementary improvements, leading to a consistent overall enhancement in MedMO's performance across all tasks.



Figure 5. **Performance across post-training stages on radiology datasets.** MedMO exhibits consistent gains in diagnostic accuracy and localization across IU-Xray, MIMIC-CXR, CheXpert, and MedTrinity datasets. The sharp improvement at Stage 2 highlights the benefit of alignment tuning with medical reasoning objectives.



Figure 6. **QA and VQA ablation across stages.** Both question-answering (QA) and visual question-answering (VQA) accuracy improve progressively, demonstrating that stage-wise optimization enhances multimodal reasoning and factual grounding in medical.

### 4.4.2. Bounding-Box Verifiable Reward

Table 4 shows consistent gains after reinforcement learning, confirming the effectiveness of our bounding-box reward. Even with small improvements, the reward reliably enhances spatial precision and grounding consistency across datasets[2].

Table 4. Absolute change ($\Delta$) after bouding box verifiable reward[2].

| Dataset | Before (IoU) | After (IoU) | $\Delta$ (IoU) |
|---|---|---|---|
| NIH | 8.8 | 13.3 | +4.5 |
| DeepLesion | 38.5 | 38.9 | +0.4 |
| Bacteria | 54.6 | 55.0 | +0.4 |

## 5. Conclusion

We introduced **MedMO**, a general-purpose medical multimodal foundation model that unifies visual grounding, clinical reasoning, and language understanding across diverse medical modalities. Through a scalable four-stage post-training pipeline, spanning large-scale alignment, high-resolution fine-tuning, instruction tuning, and reinforcement learning with verifiable rewards-MedMO achieves robust multimodal comprehension and precise spatial localization. Experimental results demonstrate substantial gains over strong open-source baselines As an open medical MLLM, MedMO establishes a scalable path toward reliable and transparent medical VLM systems. Future work could explore strategies to better retain SFT knowledge within reinforcement learning frameworks.

**Limitation.** MedMO's stage-wise training introduces minor task-level performance shifts, as shown in Figures 5 and 6, a typical behavior in large multimodal models due to *catastrophic forgetting* [55]. Future work will focus on improving cross-task retention while expanding coverage across additional medical imaging modalities.

# References

[1] Asad Aali, Dave Van Veen, YI Arefeen, Jason Hom, Christian Bluethgen, Eduardo Pontes Reis, Sergios Gatidis, Namuun Clifford, Joseph Daws, Arash Tehrani, et al. Mimic-iv-ext-bhc: labeled clinical notes dataset for hospital course summarization. *PhysioNet*, 1(0):10–13026, 2024. 6

[2] Asad Aali, Dave Van Veen, Yamin Arefeen, Jason Hom, Christian Bluethgen, Eduardo Pontes Reis, Sergios Gatidis, Namuun Clifford, Joseph Daws, Arash Tehrani, Jangwon Kim, and Akshay Chaudhari. Mimic-iv-ext-bhc: Labeled clinical notes dataset for hospital course summarization (version 1.2.0), 2025. 6

[3] Asma Ben Abacha, Sadid A Hasan, Vivek V Datla, Joey Liu, Dina Demner-Fushman, and Henning Müller. Vqa-med: Overview of the medical visual question answering task at imageclef 2019. *CLEF (working notes)*, 2(6):1–11, 2019. 6

[4] Rawan AlSaad, Alaa Abd-alrazaq, Sabri Boughorbel, Arfan Ahmed, Max-Antoine Renault, Rafat Damseh, and Javaid Sheikh. Multimodal large language models in health care: Applications, challenges, and future outlook. *J Med Internet Res*, 26, 2024. 2

[5] Rahul K. Arora, Jason Wei, Rebecca Soskin Hicks, Preston Bowman, Joaquin Quiñonero-Candela, Foivos Tsimpourlas, Michael Sharman, Meghan Shah, Andrea Vallone, Alex Beutel, Johannes Heidecke, and Karan Singhal. Healthbench: Evaluating large language models towards improved human health. *ArXiv*, abs/2505.08775, 2025. 3

[6] Omer Aydin and Enis Karaarslan. Openai chatgpt interprets radiological images: Gpt-4 as a medical doctor for a fast check-up. *ArXiv*, abs/2501.06269, 2025. 3

[7] Seongsu Bae, Daeun Kyung, Jaehee Ryu, Eunbyeol Cho, Gyubok Lee, Sunjun Kweon, Jungwoo Oh, Lei Ji, Eric Chang, Tackeun Kim, et al. Mimic-ext-mimic-cxr-vqa: A complex, diverse, and large-scale visual question answering dataset for chest x-ray images, 2024. 6

[8] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *ArXiv*, abs/2502.13923, 2025. 3

[9] Dylan Bannon, Erick Moen, Morgan Schwartz, Enrico Borba, Takamasa Kudo, Noah Greenwald, Vibha Vijayakumar, Brian Chang, Edward Pao, Erik Osterman, et al. Deepcell kiosk: scaling deep learning–enabled cellular image analysis with kubernetes. *Nature methods*, 18(1):43–45, 2021. 5

[10] Asma Ben Abacha and Dina Demner-Fushman. A question-entailment approach to question answering. *BMC Bioinform.*, 20(1):511:1–511:23, 2019. 6

[11] Harald Brodoefel, Christof Burgstahler, Ilias Tsiflikas, Anja Reimann, Stephen Schroeder, Claus D Claussen, Martin Heuschmid, and Andreas F Kopp. Dual-source ct: effect of heart rate, heart rate variability, and calcification on image quality and diagnostic accuracy. *Radiology*, 247(2):346–355, 2008. 6

[12] Pierre Chambon, Jean-Benoit Delbrouck, Thomas Sounack, Shih-Cheng Huang, Zhihong Chen, Maya Varma, Steven QH Truong, Chu The Chuong, and Curtis P Langlotz. Chexpert plus: Augmenting a large chest x-ray dataset with text radiology reports, patient demographics and additional image formats. *arXiv preprint arXiv:2405.19538*, 2024. 6

[13] Junying Chen, Zhenyang Cai, Ke Ji, Xidong Wang, Wanlong Liu, Rongsheng Wang, Jianye Hou, and Benyou Wang. Huatuogpt-o1, towards medical complex reasoning with llms, 2024. 6

[14] Junying Chen, Chi Gui, Ruyi Ouyang, Anningzhe Gao, Shunian Chen, Guiming Hardy Chen, Xidong Wang, Ruifei Zhang, Zhenyang Cai, Ke Ji, Guangjun Yu, Xiang Wan, and Benyou Wang. Huatuogpt-vision, towards injecting medical visual knowledge into multimodal llms at scale. *ArXiv*, abs/2406.19280, 2024. 2, 3

[15] Junying Chen, Zhenyang Cai, Ke Ji, Xidong Wang, Wanlong Liu, Rongsheng Wang, and Benyou Wang. Towards medical complex reasoning with LLMs through medical verifiable problems. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 14552–14573, Vienna, Austria, 2025. Association for Computational Linguistics. 3, 6

[16] Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. Januspro: Unified multimodal understanding and generation with data and model scaling. *ArXiv*, abs/2501.17811, 2025. 2

[17] Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. Generating radiology reports via memory-driven transformer. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, 2020. 6

[18] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, Lixin Gu, Xuehui Wang, Qingyun Li, Yimin Ren, Zixuan Chen, Jiapeng Luo, Jiahao Wang, Tan Jiang, Bo Wang, Conghui He, Botian Shi, Xingcheng Zhang, Han Lv, Yi Wang, Wenqi Shao, Pei Chu, Zhongying Tu, Tong He, Zhiyong Wu, Huipeng Deng, Jiaye Ge, Kai Chen, Kaipeng Zhang, Limin Wang, Min Dou, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao, Jifeng Dai, and Wenhai Wang. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *ArXiv*, abs/2412.05271, 2025. 3

[19] Matthew Christensen, Milos Vukadinovic, Neal Yuan, and David Ouyang. Vision–language foundation model for echocardiogram interpretation. *Nature Medicine*, 30: 1481–1488, 2024. 3

[20] Google DeepMind. Advancing medical ai with med-gemini, 2024. 3

[21] Google DeepMind. Gemini: A family of highly capable multimodal models. *ArXiv*, abs/2312.11805, 2025. 2, 3

[22] Google DeepMind. Gemini 2.5: Our most intelligent ai model, 2025. 2

[23] DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv*, abs/2501.12948, 2025. 2, 4

[24] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin Bransom, Kiana Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, Mark Yatskar, Chris Callison-Burch, Andrew Head, Rose Hendrix, Favyen Bastani, Eli VanderBilt, Nathan Lambert, Yvonne Chou, Arnavi Chheda, Jenna Sparks, Sam Skjonsberg, Michael Schmitz, Aaron Sarnat, Byron Bischoff, Pete Walsh, Chris Newell, Piper Wolters, Tanmay Gupta, Kuo-Hao Zeng, Jon Borchardt, Dirk Groeneveld, Crystal Nam, Sophie Lebrecht, Caitlin Wittlif, Carissa Schoenick, Oscar Michel, Ranjay Krishna, Luca Weihs, Noah A. Smith, Hannaneh Hajishirzi, Ross Girshick, Ali Farhadi, and Aniruddha Kembhavi. Molmo and pixmo: Open weights and open data for state-of-the-art vision-language models. *ArXiv*, abs/2409.17146, 2024. 3

[25] Zhuo Deng, Weihao Gao, Chucheng Chen, Zhiyuan Niu, Zheng Gong, Ruiheng Zhang, Zhenjie Cao, Fang Li, Zhaoyi Ma, Wenbin Wei, and Lan Ma. Ophglm: An ophthalmology large language-and-vision assistant. *Artif. Intell. Med.*, 157 (C), 2024. 3

[26] Ross W Filice, Anouk Stein, Carol C Wu, Veronica A Arteaga, Stephen Borstelmann, Ramya Gaddikeri, Maya Galperin-Aizenberg, Ritu R Gill, Myrna C Godoy, Stephen B Hobbs, et al. Crowdsourcing pneumothorax annotations using machine learning annotations on the nih chest x-ray dataset. *Journal of digital imaging*, 33(2):490–496, 2020. 6

[27] Jevgenij Gamper, Navid Alemi Koohbanani, Ksenija Benes, Simon Graham, Mostafa Jahanifar, Syed Ali Khurram, Ayesha Azam, Katherine Hewitt, and Nasir Rajpoot. Pannuke dataset extension, insights and baselines. *ArXiv*, abs/2003.10778, 2020. 5

[28] Lidia Garrucho, Claire-Anne Reidel, Kaisar Kushibar, Smriti Joshi, Richard Osuala, Apostolia Tsirikoglou, Maciej Bobowicz, Javier del Riego, Alessandro Catanese, Katarzyna Gwoździewicz, et al. Mama-mia: A large-scale multi-center breast cancer dce-mri benchmark dataset with expert segmentations. *arXiv e-prints*, pages arXiv–2406, 2024. 5

[29] Ibrahim Ethem Hamamci, Sezgin Er, Anjany Sekuboyina, Enis Simsar, Alperen Tezcan, Ayse Gulnihan Simsek, Sevval Nil Esirgun, Furkan Almas, Irem Doğan, Muhammed Furkan Dasdelen, et al. Generatect: Text-conditional generation of 3d chest ct volumes. In *European Conference on Computer Vision*, pages 126–143. Springer, 2024. 5

[30] Ibrahim Ethem Hamamci, Sezgin Er, Chenyu Wang, Furkan Almas, Ayse Gulnihan Simsek, Sevval Nil Esirgun, Irem Doga, Omer Faruk Durugol, Weicheng Dai, Murong Xu, Muhammed Furkan Dasdelen, Bastian Wittmann, Tamaz Amiranashvili, Enis Simsar, Mehmet Simsar, Emine Bensu Erdemir, Abdullah Alanbay, Anjany Sekuboyina, Berkan Lafci, Christian Bluethgen, Kayhan Batmanghelich, Mehmet Kemal Ozdemir, and Bjoern Menze. Developing generalist foundation models from a multimodal dataset for 3d computed tomography. *ArXiv*, abs/2403.17834, 2025. 3

[31] Xuehai He, Yichen Zhang, Luntian Mou, Eric Xing, and Pengtao Xie. Pathvqa: 30000+ questions for medical visual question answering. *ArXiv*, abs/2003.10286, 2020. 6

[32] Stephanie L. Hyland, Shruthi Bannur, Kenza Bouzid, Daniel C. Castro, Mercy Ranjit, Anton Schwaighofer, Fernando Pérez-García, Valentina Salvatelli, Shaury Srivastav, Anja Thieme, Noel Codella, Matthew P. Lungren, Maria Teodora Wetscherek, Ozan Oktay, and Javier Alvarez-Valle. Maira-1: A specialised large multimodal model for radiology report generation. *ArXiv*, abs/2311.13668, 2024. 2, 3

[33] Wisdom Oluchi Ikezogwo, Mehmet Saygin Seyfioglu, Fatemeh Ghezloo, Dylan Stefan Chan Geva, Fatwir Sheikh Mohammed, Pavan Kumar Anand, Ranjay Krishna, and Linda Shapiro. Quilt-1m: One million image-text pairs for histopathology. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. 3

[34] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, et al. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, pages 590–597, 2019. 6

[35] Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *arXiv preprint arXiv:2009.13081*, 2020. 6

[36] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. Pubmedqa: A dataset for biomedical research question answering. *ArXiv*, abs/1909.06146, 2019. 6

[37] Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Roger G Mark, and Steven Horng. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data*, 6(1):317, 2019. 6

[38] Masakata Kawai, Noriaki Ota, and Shinsuke Yamaoka. Large-scale pretraining on pathological images for fine-tuning of small pathological benchmarks. In *Workshop on Medical Image Learning with Limited and Noisy Data*, pages 257–267. Springer, 2023. 5

[39] Kimi, Angang Du, Bohong Yin, Bowei Xing, Bowen Qu, Bowen Wang, Cheng Chen, Chenlin Zhang, Chenzhuang Du, Chu Wei, Congcong Wang, Dehao Zhang, Dikang Du, Dongliang Wang, Enming Yuan, Enzhe Lu, Fang Li, Flood Sung, Guangda Wei, Guokun Lai, Han Zhu, Hao Ding, Hao Hu, Hao Yang, Hao Zhang, Haoning Wu, Haotian Yao, Haoyu Lu, Heng Wang, Hongcheng Gao, Huabin Zheng, Jiaming Li, Jianlin Su, Jianzhou Wang, Jiaqi Deng, Jiezhong Qiu, Jin Xie, Jinhong Wang, Jingyuan Liu, Junjie Yan, Kun Ouyang, Liang Chen, Lin Sui, Longhui Yu, Mengfan Dong, Mengnan Dong, Nuo Xu, Pengyu Cheng, Qizheng Gu, Runjie Zhou, Shaowei Liu, Sihan Cao, Tao Yu, Tianhui Song, Tongtong Bai, Wei Song, Weiran He, Weixiao Huang, Weixin Xu, Xiaokun Yuan, Xingcheng Yao, Xingzhe Wu, Xinxing Zu, Xinyu Zhou, Xinyuan Wang, Y. Charles, Yan Zhong, Yang Li, Yangyang Hu, Yanru Chen,

Yejie Wang, Yibo Liu, Yibo Miao, Yidao Qin, Yimin Chen, Yiping Bao, Yiqin Wang, Yongsheng Kang, Yuanxin Liu, Yulun Du, Yuxin Wu, Yuzhi Wang, Yuzi Yan, Zaida Zhou, Zhaowei Li, Zhejun Jiang, Zheng Zhang, Zhilin Yang, Zhiqi Huang, Zihao Huang, Zijia Zhao, Ziwei Chen, and Zongyu Lin. Kimi-vl technical report. *ArXiv*, abs/2504.07491, 2025. 2

[40] Yuxiang Lai, Jike Zhong, Ming Li, Shitian Zhao, and Xiaofeng Yang. Med-r1: Reinforcement learning for generalizable medical reasoning in vision-language models. *ArXiv*, abs/2503.13939, 2025. 2, 3

[41] Jason J Lau, Soumya Gayen, Asma Ben Abacha, and Dina Demner-Fushman. A dataset of clinically generated visual questions and answers about radiology images. *Scientific data*, 5(1):1–10, 2018. 6

[42] Lavita AI. Chatdoctor-healthcaremagic-100k. https://huggingface.co/datasets/lavita/ ChatDoctor-HealthCareMagic-100k, 2023. Accessed: 2025-11-17. 6

[43] Peter Lee, Sebastien Bubeck, and Joseph Petro. Benefits, limits, and risks of gpt-4 as an ai chatbot for medicine. *New England Journal of Medicine*, 388(13):1233–1239, 2023. 2

[44] Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, Joe Davison, Mario Šaško, Gunjan Chhablani, Bhavitvya Malik, Simon Brandeis, Teven Le Scao, Victor Sanh, Canwen Xu, Nicolas Patry, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger, Clément Delangue, Théo Matussière, Lysandre Debut, Stas Bekman, Pierric Cistac, Thibault Goehringer, Victor Mustar, François Lagunas, Alexander Rush, and Thomas Wolf. Datasets: A community library for natural language processing. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 175–184, Online and Punta Cana, Dominican Republic, 2021. Association for Computational Linguistics. 6

[45] Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. LLaVA-med: Training a large language-and-vision assistant for biomedicine in one day. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. 2, 3

[46] Tianbin Li, Yanzhou Su, Wei Li, Bin Fu, Zhe Chen, Ziyan Huang, Guoan Wang, Chenglong Ma, Ying Chen, Ming Hu, Yanjun Li, Pengcheng Chen, Xiaowei Hu, Zhongying Deng, Yuanfeng Ji, Jin Ye, Yu Qiao, and Junjun He. Gmai-vl & gmai-vl-5.5m: A large vision-language model and a comprehensive multimodal dataset towards general medical ai. *ArXiv*, abs/2411.14522, 2025. 2, 3

[47] Tianwei Lin, Wenqiao Zhang, Sijing Li, Yuqian Yuan, Binhe Yu, Haoyuan Li, Wanggui He, Hao Jiang, Mengze Li, Xiaohui Song, Siliang Tang, Jun Xiao, Hui Lin, Yueting Zhuang, and Beng Chin Ooi. Healthgpt: A medical large vision-language model for unifying comprehension and generation via heterogeneous knowledge adaptation. *ArXiv*, abs/2502.09838, 2025. 2, 3

[48] Weixiong Lin, Ziheng Zhao, Xiaoman Zhang, Chaoyi Wu, Ya Zhang, Yanfeng Wang, and Weidi Xie. Pmc-clip: Contrastive language-image pre-training using biomedical documents. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 525–536. Springer, 2023. 6

[49] Bo Liu, Li-Ming Zhan, Li Xu, Lin Ma, Yan Yang, and Xiao-Ming Wu. Slake: A semantically-labeled knowledge-enhanced dataset for medical visual question answering. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pages 1650–1654, 2021. 6

[50] Junling Liu, Ziming Wang, Qichen Ye, Dading Chong, Peilin Zhou, and Yining Hua. Qilin-med-vl: Towards chinese large vision-language model for general healthcare. *ArXiv*, abs/2310.17956, 2023. 2

[51] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European conference on computer vision*, pages 38–55. Springer, 2024. 3

[52] Meng Lou, Hanning Ying, Xiaoqing Liu, Hong-Yu Zhou, Yuqin Zhang, and Yizhou Yu. Sdr-former: A siamese dual-resolution transformer for liver lesion classification using 3d multi-phase imaging. *Neural Networks*, page 107228, 2025. 5

[53] Ming Y. Lu, Bowen Chen, Drew F. K. Williamson, Richard J. Chen, Kenji Ikamura, Georg Gerber, Ivy Liang, Long Phi Le, Tong Ding, Anil V Parwani, and Faisal Mahmood. A multimodal generative ai copilot for human pathology. *Nature*, 634:466–473, 2024. 3

[54] Yan Luo, Min Shi, Muhammad Osama Khan, Muhammad Muneeb Afzal, Hao Huang, Shuaihang Yuan, Yu Tian, Luo Song, Ava Kouhana, Tobias Elze, et al. Fairclip: Harnessing fairness in vision-language learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12289–12301, 2024. 6

[55] Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. An empirical study of catastrophic forgetting in large language models during continual fine-tuning. *IEEE Transactions on Audio, Speech and Language Processing*, 2025. 8

[56] MedAlpaca. medical_meadow_medical_flashcards. https://huggingface.co/datasets/ medalpaca/medical_meadow_medical_ flashcards, 2023. Accessed: 2025-11-17. 6

[57] MedAlpaca. medical_meadow_wikidoc. https: //huggingface.co/datasets/medalpaca/ medical_meadow_wikidoc, 2023. Accessed: 2025-11-17. 6

[58] Michael Moor, Qian Huang, Shirley Wu, Michihiro Yasunaga, Yash Dalmia, Jure Leskovec, Cyril Zakka, Eduardo Pontes Reis, and Pranav Rajpurkar. Med-flamingo: a multimodal medical few-shot learner. In *Proceedings of the 3rd Machine Learning for Health Symposium*, pages 353–367. PMLR, 2023. 2

[59] Eszter Nagy, Michael Janisch, Franko Hržić, Erich Sorantin, and Sebastian Tschauner. A pediatric wrist trauma x-ray

dataset (grazpedwri-dx) for machine learning. *Scientific data*, 9(1):222, 2022. 6

[60] Vishwesh Nath, Wenqi Li, Dong Yang, Andriy Myronenko, Mingxin Zheng, Yao Lu, Zhijian Liu, Hongxu Yin, Yucheng Tang, Pengfei Guo, Can Zhao, Ziyue Xu, Yufan He, Greg Heinrich, Yee Man Law, Benjamin Simon, Stephanie Harmon, Stephen Aylward, Marc Edgar, Michael Zephyr, Song Han, Pavlo Molchanov, Baris Turkbey, Holger Roth, and Daguang Xu. Vila-m3: Enhancing vision-language models with medical expert knowledge. *ArXiv*, abs/2411.12915, 2025. 3

[61] Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. Capabilities of gpt-4 on medical challenge problems. *ArXiv*, abs/2303.13375, 2023. 2

[62] OpenAI. Gpt-4o system card. *ArXiv*, abs/2410.21276, 2024. 2

[63] OpenAI. Openai o1 system card. *ArXiv*, abs/2412.16720, 2024. 2

[64] OpenAI. Introducing gpt-4.1 in the api, 2025.

[65] OpenAI. Introducing *o3* and *o4-mini*, 2025. 2

[66] Suraj Pai, Ibrahim Hadzic, Dennis Bontempi, Keno Bressem, Benjamin H. Kann, Andriy Fedorov, Raymond H. Mak, and Hugo J. W. L. Aerts. Vision foundation models for computed tomography. *ArXiv*, abs/2501.09001, 2025. 3

[67] Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Proceedings of the Conference on Health, Inference, and Learning*, pages 248–260. PMLR, 2022. 6

[68] Jiazhen Pan, Che Liu, Junde Wu, Fenglin Liu, Jiayuan Zhu, Hongwei Bran Li, Chen Chen, Cheng Ouyang, and Daniel Rueckert. Medvlm-r1: Incentivizing medical reasoning capability of vision-language models (vlms) via reinforcement learning. *ArXiv*, abs/2502.19634, 2025. 2, 3

[69] Obioma Pelka, Sven Koitka, Johannes Rückert, Felix Nensa, and Christoph M. Friedrich. Radiology objects in context (roco): A multimodal image dataset. In *Intravascular Imaging and Computer Assisted Stenting and Large-Scale Annotation of Biomedical Data and Expert Label Synthesis*, pages 180–189, Cham, 2018. Springer International Publishing. 6

[70] Raymond D Ratliff. Meadows in the sierra nevada of california: state of knowledge. 1985. 6

[71] Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 658–666, 2019. 5, 1

[72] Johannes Rückert, Louise Bloch, Raphael Brüngel, Ahmad Idrissi-Yaghir, Henning Schäfer, Cynthia S. Schmidt, Sven Koitka, Obioma Pelka, Asma Ben Abacha, Alba G. Seco de Herrera, Henning Müller, Peter A. Horn, Felix Nensa, and Christoph M. Friedrich. Rocov2: Radiology objects in context version 2, an updated multimodal image dataset. *Scientific Data*, 11(1), 2024. 6

[73] Khaled Saab, Tao Tu, Wei-Hung Weng, Ryutaro Tanno, David Stutz, Ellery Wulczyn, Fan Zhang, Tim Strother, Chunjong Park, Elahe Vedadi, Juanma Zambrano Chaves, Szu-Yeu Hu, Mike Schaekermann, Aishwarya Kamath, Yong Cheng, David G. T. Barrett, Cathy Cheung, Basil Mustafa, Anil Palepu, Daniel McDuff, Le Hou, Tomer Golany, Luyang Liu, Jean baptiste Alayrac, Neil Houlsby, Nenad Tomasev, Jan Freyberg, Charles Lau, Jonas Kemp, Jeremy Lai, Shekoofeh Azizi, Kimberly Kanada, SiWai Man, Kavita Kulkarni, Ruoxi Sun, Siamak Shakeri, Luheng He, Ben Caine, Albert Webson, Natasha Latysheva, Melvin Johnson, Philip Mansfield, Jian Lu, Ehud Rivlin, Jesper Anderson, Bradley Green, Renee Wong, Jonathan Krause, Jonathon Shlens, Ewa Dominowska, S. M. Ali Eslami, Katherine Chou, Claire Cui, Oriol Vinyals, Koray Kavukcuoglu, James Manyika, Jeff Dean, Demis Hassabis, Yossi Matias, Dale Webster, Joelle Barral, Greg Corrado, Christopher Semturs, S. Sara Mahdavi, Juraj Gottweis, Alan Karthikesalingam, and Vivek Natarajan. Capabilities of gemini models in medicine. *ArXiv*, abs/2404.18416, 2024. 3

[74] Mourad Sarrouti. Nlm at vqa-med 2020: Visual question answering and generation in the medical domain. In *CLEF (Working Notes)*, 2020. 6

[75] Mehmet Saygin Seyfioglu, Wisdom O Ikezogwo, Fatemeh Ghezloo, Ranjay Krishna, and Linda Shapiro. Quilt-llava: Visual instruction tuning by extracting localized narratives from open-source histopathology videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13183–13192, 2024. 6

[76] Mehmet Saygin Seyfioglu, Wisdom O. Ikezogwo, Fatemeh Ghezloo, Ranjay Krishna, and Linda Shapiro. Quilt-llava: Visual instruction tuning by extracting localized narratives from open-source histopathology videos. *ArXiv*, abs/2312.04746, 2025. 2, 3

[77] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *ArXiv*, abs/2402.03300, 2024. 2, 4

[78] Yan Shu, Chi Liu, Robin Chen, Derek Li, and Bryan Dai. Fleming-vl: Towards universal medical visual reasoning with multimodal llms. *arXiv preprint arXiv:2511.00916*, 2025. 3

[79] Zhongyi Shui, Zhongyi Shui, Jianpeng Zhang, Weiwei Cao, Sinuo Wang, Ruizhe Guo, Le Lu, Lin Yang, Xianghua Ye, Tingbo Liang, Qi Zhang, and Ling Zhang. Large-scale and fine-grained vision-language pre-training for enhanced ct image understanding. In *The Thirteenth International Conference on Learning Representations*, 2025. 3

[80] Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Abubakr Babiker, Nathanael Schärli, Aakanksha Chowdhery, Philip Mansfield, Dina Demner-Fushman, Blaise Agüera y Arcas, Dale Webster, Greg S. Corrado, Yossi Matias, Katherine Chou, Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Rajkomar, Joelle Barral, Christopher Semturs, Alan Karthikesalingam, and

Vivek Natarajan. Large language models encode clinical knowledge. *Nature*, 620:172–180, 2023. 3

[81] Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, Mike Schaekermann, Amy Wang, Mohamed Amin, Sami Lachgar, Philip Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska, Blaise Aguera y Arcas, Nenad Tomasev, Yun Liu, Renee Wong, Christopher Semturs, S. Sara Mahdavi, Joelle Barral, Dale Webster, Greg S. Corrado, Yossi Matias, Shekoofeh Azizi, Alan Karthikesalingam, and Vivek Natarajan. Towards expert-level medical question answering with large language models. *Nature Medicine*, 31:943–950, 2023. 3

[82] I Siragusa, S Contino, ML Ciura, R Alicata, and R MedPix Pirrone. 2.0: A comprehensive multimodal biomedical data set for advanced ai applications. arxiv 2024. *arXiv preprint arXiv:2407.02994*. 6

[83] Ryutaro Tanno, David G. T. Barrett, Andrew Sellergren, Sumedh Ghaisas, Sumanth Dathathri, Abigail See, Johannes Welbl, Charles Lau, Tao Tu, Shekoofeh Azizi, Karan Singhal, Mike Schaekermann, Rhys May, Roy Lee, SiWai Man, Sara Mahdavi, Zahra Ahmed, Yossi Matias, Joelle Barral, S. M. Ali Eslami, Danielle Belgrave, Yun Liu, Sreenivasa Raju Kalidindi, Shravya Shetty, Vivek Natarajan, Pushmeet Kohli, Po-Sen Huang, Alan Karthikesalingam, and Ira Ktena. Collaboration between clinicians and vision–language models in radiology report generation. *Nature Medicine*, 31, 2025. 3

[84] Dianzhe Tian, Shitao Jiang, Lei Zhang, Xin Lu, and Yiyao Xu. The role of large language models in medical image processing: a narrative review. *Quantitative Imaging in Medicine and Surgery*, 14(1):1108, 2023. 2

[85] Yuri Tolkach, Lisa Marie Wolgast, Alexander Damanakis, Alexey Pryalukhin, Simon Schallenberg, Wolfgang Hulla, Marie-Lisa Eich, Wolfgang Schroeder, Anirban Mukhopadhyay, Moritz Fuchs, et al. Artificial intelligence for tumour tissue detection and histological regression grading in oesophageal adenocarcinomas: a retrospective algorithm development and validation study. *The Lancet Digital Health*, 5(5):e265–e275, 2023. 5

[86] Katarzyna Tomczak, Patrycja Czerwińska, and Maciej Wiznerowicz. Review the cancer genome atlas (tcga): an immeasurable source of knowledge. *Contemporary Oncology/Współczesna Onkologia*, 2015(1):68–77, 2015. 5

[87] Simon van Vliet, Annina R Winkler, Stefanie Spriewald, Bärbel Stecher, Martin Ackermann, et al. Spatially correlated gene expression in bacterial groups: the role of lineage history, spatial gradients, and cell-cell interactions. *Cell systems*, 6(4):496–507, 2018. 5, 6

[88] Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, Shengyi Huang, Kashif Rasul, and Quentin Gallouédec. Trl: Transformer reinforcement learning. https://github.com/huggingface/trl, 2020. 5

[89] Patrick Wagner, Maximilian Springenberg, Marius Kröger, Rose KC Moritz, Johannes Schleusener, Martina C Meinke, and Jackie Ma. Semantic modeling of cell damage prediction: a machine learning approach at human-level perfor-

mance in dermatology. *Scientific Reports*, 13(1):8336, 2023. 5

[90] Lucy Lu Wang, Kyle Lo, Yoganand Chandrasekhar, Russell Reas, Jiangjiang Yang, Doug Burdick, Darrin Eide, Kathryn Funk, Yannis Katsis, Rodney Michael Kinney, Yunyao Li, Ziyang Liu, William Merrill, Paul Mooney, Dewey A. Murdick, Devvret Rishi, Jerry Sheehan, Zhihong Shen, Brandon Stilson, Alex D. Wade, Kuansan Wang, Nancy Xin Ru Wang, Christopher Wilhelm, Boya Xie, Douglas M. Raymond, Daniel S. Weld, Oren Etzioni, and Sebastian Kohlmeier. CORD-19: The COVID-19 open research dataset. In *Proceedings of the 1st Workshop on NLP for COVID-19 at ACL 2020*, Online, 2020. Association for Computational Linguistics. 6

[91] Sheng Wang, Zihao Zhao, Xi Ouyang, Tianming Liu, Qian Wang, and Dinggang Shen. Interactive computer-aided diagnosis on medical image using large language models. *Communications Engineering*, 3, 2024. 3

[92] Xiyue Wang, Junhan Zhao, Eliana Marostica, Wei Yuan, Jietian Jin, Jiayu Zhang, Ruijiang Li, Hongping Tang, Kanran Wang, Yu Li, Fang Wang, Yulong Peng, Junyou Zhu, Jing Zhang, Christopher R. Jackson, Jun Zhang, Deborah Dillon, Nancy U. Lin, Lynette Sholl, Thomas Denize, David Meredith, Keith L. Ligon, Sabina Signoretti, Shuji Ogino, Jeffrey A. Golden, MacLean P. Nasrallah, Xiao Han, Sen Yang, and Kun-Hsing Yu. A pathology foundation model for cancer diagnosis and prognosis prediction. *Nature*, 634: 970–978, 2024. 3

[93] Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Towards generalist foundation model for radiology by leveraging web-scale 2d&3d medical data. *ArXiv*, abs/2308.02463, 2023. 2

[94] Juncheng Wu, Wenlong Deng, Xingxuan Li, Sheng Liu, Taomian Mi, Yifan Peng, Ziyang Xu, Yi Liu, Hyunjin Cho, Chang-In Choi, Yihan Cao, Hui Ren, Xiang Li, Xiaoxiao Li, and Yuyin Zhou. Medreason: Eliciting factual medical reasoning steps in llms via knowledge graphs. *ArXiv*, abs/2504.00993, 2025. 6

[95] Yunfei Xie, Juncheng Wu, Haoqin Tu, Siwei Yang, Bingchen Zhao, Yongshuo Zong, Qiao Jin, Cihang Xie, and Yuyin Zhou. A preliminary study of o1 in medicine: Are we closer to an ai doctor? *ArXiv*, abs/2409.15277, 2024. 3

[96] Yunfei Xie, Ce Zhou, Lang Gao, Juncheng Wu, Xianhang Li, Hong-Yu Zhou, Sheng Liu, Lei Xing, James Zou, Cihang Xie, and Yuyin Zhou. Medtrinity-25m: A large-scale multimodal dataset with multigranular annotations for medicine. In *The Thirteenth International Conference on Learning Representations*, 2025. 3, 5, 6

[97] Weiwen Xu, Hou Pong Chan, Long Li, Mahani Aljunied, Ruifeng Yuan, Jianyu Wang, Chenghao Xiao, Guizhen Chen, Chaoqun Liu, Zhaodonghui Li, et al. Lingshu: A generalist foundation model for unified multimodal medical understanding and reasoning. *arXiv preprint arXiv:2506.07044*, 2025. 3

[98] Ke Yan, Xiaosong Wang, Le Lu, and Ronald M Summers. Deeplesion: automated mining of large-scale lesion annotations and universal lesion detection with deep learning. *Journal of medical imaging*, 5(3):036501–036501, 2018. 6

[99] Zhiling Yan, Kai Zhang, Rong Zhou, Lifang He, Xiang Li, and Lichao Sun. Multimodal chatgpt for medical applications: an experimental study of gpt-4v. *ArXiv*, abs/2310.19061, 2023. 2

[100] Hang Yang, Hao Chen, Hui Guo, Yineng Chen, Ching-Sheng Lin, Shu Hu, Jinrong Hu, Xi Wu, and Xin Wang. Llm-medqa: Enhancing medical question answering through case studies in large language models. *arXiv preprint arXiv:2501.05464*, 2024. 6

[101] Lin Yang, Shawn Xu, Andrew Sellergren, Timo Kohlberger, Yuchen Zhou, Ira Ktena, Atilla Kiraly, Faruk Ahmed, Farhad Hormozdiari, Tiam Jaroensri, Eric Wang, Ellery Wulczyn, Fayaz Jamil, Theo Guidroz, Chuck Lau, Siyuan Qiao, Yun Liu, Akshay Goel, Kendall Park, Arnav Agharwal, Nick George, Yang Wang, Ryutaro Tanno, David G. T. Barrett, Wei-Hung Weng, S. Sara Mahdavi, Khaled Saab, Tao Tu, Sreenivasa Raju Kalidindi, Mozziyar Etemadi, Jorge Cuadros, Gregory Sorensen, Yossi Matias, Katherine Chou, Greg Corrado, Joelle Barral, Shravya Shetty, David Fleet, S. M. Ali Eslami, Daniel Tse, Shruthi Prabhakara, Cory McLean, Dave Steiner, Rory Pilgrim, Christopher Kelly, Shekoofeh Azizi, and Daniel Golden. Advancing multimodal medical capabilities of gemini. *arXiv*, abs/2405.03162, 2024. 3

[102] YongchengYAO. Kipa22. https://huggingface.co/datasets/YongchengYAO/KiPA22, 2025. Accessed: 2025-11-17. 5

[103] Bei Yu, Yingya Li, and Jun Wang. Detecting causal language use in science findings. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4664–4674, Hong Kong, China, 2019. Association for Computational Linguistics. 6

[104] Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Jinhua Zhu, Jiaze Chen, Jiangjie Chen, Chengyi Wang, Hongli Yu, Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei-Ying Ma, Ya-Qin Zhang, Lin Yan, Mu Qiao, Yonghui Wu, and Mingxuan Wang. Dapo: An open-source llm reinforcement learning system at scale. *ArXiv*, abs/2503.14476, 2025. 2, 4

[105] Jingkun Yue, Siqi Zhang, Zinan Jia, Huihuan Xu, Zongbo Han, Xiaohong Liu, and Guangyu Wang. Medsg-bench: A benchmark for medical image sequences grounding. *arXiv preprint arXiv:2505.11852*, 2025. 3, 6

[106] Juan Manuel Zambrano Chaves, Shih-Cheng Huang, Yanbo Xu, Hanwen Xu, Naoto Usuyama, Sheng Zhang, Fei Wang, Yujia Xie, Mahmoud Khademi, Ziyi Yang, Hany Awadalla, Julia Gong, Houdong Hu, Jianwei Yang, Chunyuan Li, Jianfeng Gao, Yu Gu, Cliff Wong, Mu Wei, Tristan Naumann, Muhao Chen, Matthew P. Lungren, Akshay Chaudhari, Serena Yeung-Levy, Curtis P. Langlotz, Sheng Wang, and Hoifung Poon. A clinically accessible small multimodal radiology model and evaluation metric for chest x-ray findings. *Nature Communications*, 16(1), 2025. 3

[107] Kai Zhang, Rong Zhou, Eashan Adhikarla, Zhiling Yan, Yixin Liu, Jun Yu, Zhengliang Liu, Xun Chen, Brian D. Davison, Hui Ren, Jing Huang, Chen Chen, Yuyin Zhou, Sunyang Fu, Wei Liu, Tianming Liu, Xiang Li, Yong Chen, Lifang He, James Zou, Quanzheng Li, Hongfang Liu, and Lichao Sun. A generalist vision–language foundation model for diverse biomedical tasks. *Nature Medicine*, 30(11):3129–3141, 2024. 2

[108] Xinlu Zhang, Chenxin Tian, Xianjun Yang, Lichang Chen, Zekun Li, and Linda Ruth Petzold. Alpacare: Instruction-tuned large language models for medical application, 2023. 6

[109] Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Weixiong Lin, Ya Zhang, Yanfeng Wang, and Weidi Xie. Pmc-vqa: Visual instruction tuning for medical visual question answering. *arXiv preprint arXiv:2305.10415*, 2023. 6

[110] Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Weixiong Lin, Ya Zhang, Yanfeng Wang, and Weidi Xie. Pmc-vqa: Visual instruction tuning for medical visual question answering. *ArXiv*, abs/2305.10415, 2024. 2

[111] Theodore Zhao, Yu Gu, Jianwei Yang, Naoto Usuyama, Ho Hin Lee, Sid Kiblawi, Tristan Naumann, Jianfeng Gao, Angela Crabtree, Jacob Abel, Christine Moung-Wen, Brian Piening, Carlo Bifulco, Mu Wei, Hoifung Poon, and Sheng Wang. A foundation model for joint segmentation, detection and recognition of biomedical objects across nine modalities. *Nature Methods*, 22(1):166–176, 2024. 3

[112] Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, Zhangwei Gao, Erfei Cui, Xuehui Wang, Yue Cao, Yangzhou Liu, Xingguang Wei, Hongjie Zhang, Haomin Wang, Weiye Xu, Hao Li, Jiahao Wang, Nianchen Deng, Songze Li, Yinan He, Tan Jiang, Jiapeng Luo, Yi Wang, Conghui He, Botian Shi, Xingcheng Zhang, Wenqi Shao, Junjun He, Yingtong Xiong, Wenwen Qu, Peng Sun, Penglong Jiao, Han Lv, Lijun Wu, Kaipeng Zhang, Huipeng Deng, Jiaye Ge, Kai Chen, Limin Wang, Min Dou, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao, Jifeng Dai, and Wenhai Wang. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *ArXiv*, abs/2504.10479, 2025. 2

# Appendix

## A. Reward function details

### A.1. Bounding Box Reward Function

For grounding tasks in the reinforcement learning stage, we employ a specialized reward function that evaluates the quality of predicted bounding boxes against ground truth annotations. This reward is computed using Hungarian matching combined with geometric metrics.

**Notation and Setup.** Given ground truth boxes $\mathcal{G} = \{g_j\}_{j=1}^{G}$ and predicted boxes $\mathcal{P} = \{p_i\}_{i=1}^{P}$ in XYXY format (i.e., $(x_1, y_1, x_2, y_2)$ coordinates), we first determine the image dimensions $(H, W)$ from the maximum extents of ground truth boxes if available, otherwise from predictions (with fallback to $(1, 1)$ if both are empty).

**Pairwise Metrics.** For each pair of boxes $(p_i, g_j)$, we compute two geometric measures:

**Normalized L1 Distance:** The L1 distance over all four coordinates, normalized by the image perimeter:

$$L1_{ij} = \frac{|x_1^p - x_1^g| + |y_1^p - y_1^g| + |x_2^p - x_2^g| + |y_2^p - y_2^g|}{2\sqrt{H^2 + W^2}} \tag{7}$$

**Generalized IoU (GIoU):** We compute $\text{GIoU}_{ij} \in [-1, 1]$ following Rezatofighi et al. [71], which extends standard IoU to account for non-overlapping boxes.

**Hungarian Matching.** To establish optimal correspondence between predictions and ground truth, we construct a cost matrix:

$$C_{ij} = w_{\text{L1}}^m \cdot L1_{ij} + w_{\text{G}}^m \cdot (1 - \text{GIoU}_{ij}), \tag{8}$$

where $w_{\text{L1}}^m = 5.0$ and $w_{\text{G}}^m = 2.0$ are matching cost weights. We apply the Hungarian algorithm to find the minimum-cost bipartite matching, yielding $m = \min(P, G)$ matched pairs $\{(i_k, j_k)\}_{k=1}^{m}$.

**Per-Match Score.** For each matched pair $(i_k, j_k)$, we compute a quality score by:

1. Mapping GIoU to $[0, 1]$: $\tilde{G}_k = \frac{\text{GIoU}_{i_k j_k} + 1}{2}$
2. Clamping L1 to $[0, 1]$: $\hat{L1}_k = \text{clip}_{[0,1]}(L1_{i_k j_k})$
3. Computing weighted blend:

$$s_k = \frac{w_{\text{L1}} \cdot (1 - \hat{L1}_k) + w_{\text{G}} \cdot \tilde{G}_k}{w_{\text{L1}} + w_{\text{G}}}, \quad s_k \in [0, 1] \tag{9}$$

where $w_{\text{L1}} = 5.0$ and $w_{\text{G}} = 2.0$ are pair score weights.

**Final Reward Computation.** The base reward is the coverage-normalized sum of matched pair scores:

$$\text{base} = \frac{1}{G} \sum_{k=1}^{m} s_k \tag{10}$$

We optionally apply penalties for false positives (FP) and false negatives (FN):

$$\text{penalty} = \frac{\lambda_{\text{FN}} \cdot (G - m) + \lambda_{\text{FP}} \cdot (P - m)}{\max(1, G)}, \tag{11}$$

where $\lambda_{\text{FN}}$ and $\lambda_{\text{FP}}$ are penalty coefficients (default: 0). The final bounding box reward is:

$$\boxed{R_{\text{bbox}} = \text{clip}_{[0,1]}(\text{base} - \text{penalty})} \tag{12}$$

Expanding the base term:

$$\text{base} = \frac{1}{G} \sum_{k=1}^{m} \frac{w_{\text{L1}}(1 - L1_{i_k j_k}) + w_{\text{G}} \cdot \frac{\text{GIoU}_{i_k j_k} + 1}{2}}{w_{\text{L1}} + w_{\text{G}}} \tag{13}$$

**Edge Cases.** The reward function handles special cases as follows:
- **No ground truth boxes** ($G = 0$): $R_{\text{bbox}} = 0.5$ (neutral reward)
- **Ground truth present but no predictions** ($G > 0, P = 0$): $R_{\text{bbox}} = \text{clip}_{[0,1]}(0 - \text{penalty})$, which equals $0.0$ with default penalties
- **Failed matching** (no feasible pairs): Treated as $m = 0$, where all ground truth boxes are unmatched and all predictions are false positives

This reward formulation encourages the model to produce accurate bounding box predictions through Hungarian-matched optimization of both localization (L1) and overlap quality (GIoU), while penalizing missing detections and spurious predictions.

## B. Experimental Details

We conducted all experiments using the SFT_Trainer and RL (GRPO) trainer frameworks. Unless otherwise noted, we used mixed-precision training (dtype=bfloat16) on a cluster of $64\times$ AMD Instinct MI210 GPUs. Random seeds, optimizer state, and scheduler configuration were logged for full reproducibility.

### B.1. Stage 1: General SFT

**Parameters Details**

We provide detailed experimental settings in Table 5, which we apply exclusively to training stage 1 MedMO.

| Parameter | Value |
| --- | --- |
| Batch size | 10 |
| Gradient accumulation steps | 2 |
| Learning rate (initial) | $1 \times 10^{-5}$ |
| LR scheduler | Cosine decay |
| Number of epochs | 1 |
| Image resolution | $768 \times 768$ pixels |
| dtype | bfloat16 |

Table 5. Training parameter details for stage 1.

## Training Dynamics

During Stage 1, optimization converges rapidly: the loss drops from $\sim 11$ to $< 0.3$ within the first $\approx 10$ steps, and entropy collapses from $\sim 5.3$ to $\sim 0.1$ over the same window, indicating quickly sharpened token distributions. Mean token accuracy rises steeply from $\sim 0.6$ to $\sim 0.95$ by step $\approx 10$ and then plateaus with minor oscillations thereafter. These curves reflect stable optimization under the cosine schedule, fast fit to the instruction format, and no signs of late-stage instability during the single-epoch SFT. *Unless noted, one plotted "step" corresponds to an aggregate over 100 mini-batches (logging interval = 100 batches).*
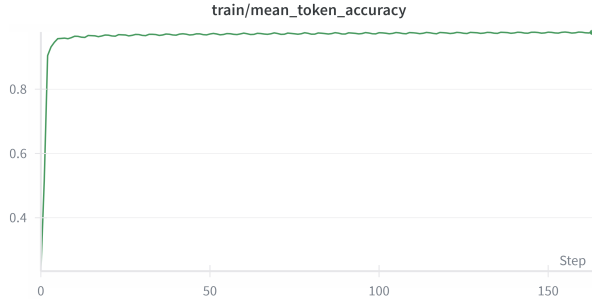


Figure 7. Stage 1 mean token accuracy vs. step (each step = 100 mini-batches). Accuracy jumps to $\sim 0.95$ within $\approx 10$ steps and remains stable.

## B.2. Stage 2: High-Resolution Image SFT

### Parameters Details

We provide detailed experimental settings in Table 6, which we apply exclusively to training stage 2 MedMO.

### Training Dynamics

During Stage 2, we fine-tuned MedMO on high-resolution ($1280 \times 1280$) medical images using a combination of VQA, grounding, and report-generation datasets. Each logged step corresponds to 100 training batches. As illustrated in Figures 10–12, the model exhibits rapid convergence and stable learning behavior. Mean token accuracy (Fig. 10) increases



Figure 8. Stage 1 training loss vs. step (each step = 100 mini-batches). Loss declines from $\sim 11$ to $< 0.3$ in the first $\approx 10$ steps, then flattens.



Figure 9. Stage 1 output entropy vs. step (each step = 100 mini-batches). Entropy collapses from $\sim 5.3$ to $\sim 0.1$ by $\approx 10$ steps, indicating confident token distributions.

| Parameter | Value |
| --- | --- |
| Batch size | 2 |
| Gradient accumulation steps | 8 |
| Learning rate (initial) | $8 \times 10^{-6}$ |
| LR scheduler | Cosine decay |
| Number of epochs | 1 |
| Image resolution | $1280 \times 1280$ pixels |
| dtype | bfloat16 |

Table 6. Training parameter details for stage 2.

sharply from $\sim 0.86$ to $\sim 0.95$ within the first few hundred steps, indicating strong adaptation to high-resolution visual–textual data. Training loss (Fig. 11) decreases quickly from $\sim 0.9$ to $\sim 0.3$ and then plateaus, confirming smooth optimization without overfitting. Entropy (Fig. 12) drops from $\sim 0.65$ to $\sim 0.27$ and remains steady, showing reduced uncertainty and confident token predictions. These results confirm that Stage 2 effectively enhances MedMO's multimodal alignment and high-resolution spatial reasoning.

Figure 10. Stage 2 mean token accuracy vs. global step (each step = 100 mini-batches). Accuracy improves rapidly from ∼0.86 to ∼0.95, showing strong convergence and model stability.



Figure 11. Stage 2 training loss vs. global step (each step = 100 mini-batches). Loss decreases from ∼0.9 to ∼0.3, confirming efficient optimization and stable convergence.
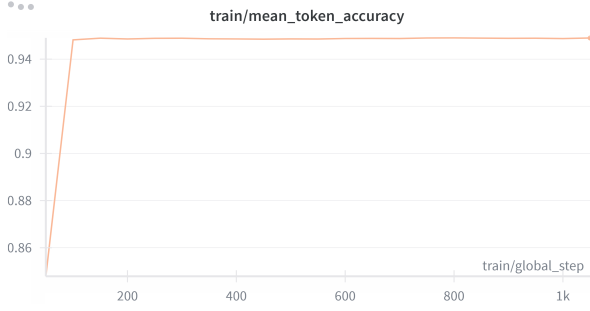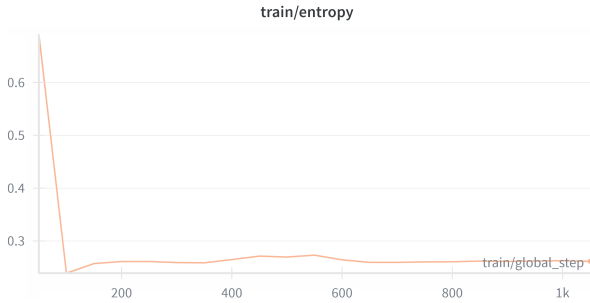


Figure 12. Stage 2 output entropy vs. global step (each step = 100 mini-batches). Entropy declines from ∼0.65 to ∼0.27, reflecting reduced uncertainty and higher confidence in predictions.

**Datasets Used**

For Stage 2, we employed datasets emphasizing multimodal reasoning, high-quality medical captions, and spatial grounding. The training corpus included a diverse mix of **VQA-oriented datasets** such as *VQA-Med-2019*, *PubMed-Vision*, *NIH-VQA*, *Quilt-LLaVA-Pretrain*, *MIMIC-Ext-MIMIC-CXR-VQA*, *VQA-RAD*, *PathVQA*, *PMC-VQA*, *SLAKE*, and *CT-*

*RATE*. We also incorporated **report-generation datasets** including *IU-Xray*, *MIMIC-CXR*, *CheXpert*, *CheXpert Plus*, *MEDPIX-ClinQA*, *ROCO*, *ROCO-V2*, and *FairVLMed* to enhance radiology-style narrative generation and image–text consistency. Finally, for **grounding and bounding-box prediction**, we used *NIH Chest X-ray*, *DeepLesion*, *GRAZPEDWRI-DX*, *SLAKE*, *Cell Microscopy (DeepCell, Bacteria, and CTC)*, and *MedSG*, which provide localized annotations for spatial reasoning and fine-grained object detection.

This combination allows MedMO to improve fine-grained visual grounding and detailed report synthesis under high-resolution supervision.

### B.3. Stage 3: Instruction Tuning

**Parameters Details**

We provide detailed experimental settings in Table 7, which we apply exclusively to training stage 3 MedMO.

| Parameter | Value |
| --- | --- |
| Batch size | 14 |
| Gradient accumulation steps | 2 |
| Learning rate (initial) | $5 \times 10^{-6}$ |
| LR scheduler | Cosine decay |
| Number of epochs | 1 |
| dtype | bfloat16 |

Table 7. Training parameter details for stage 2.

**Training Dynamics**

Stage 3 focuses on instruction tuning to enhance MedMO's clinical reasoning, comprehension, and text generation capabilities. Each step shown in the plots corresponds to 100 mini-batches. As shown in Figures 13–15, the model exhibits smooth and stable convergence. Mean token accuracy (Fig. 13) rises steadily from ∼0.62 to ∼0.69, demonstrating improved instruction-following and cross-modal reasoning. Training loss (Fig. 14) decreases from ∼1.7 to ∼1.4 within the first few steps, while entropy (Fig. 15) declines from ∼1.55 to ∼1.38, both indicating effective optimization and improved confidence. Overall, Stage 3 consolidates multimodal understanding and instruction-following capabilities with stable convergence and balanced learning dynamics.

**Datasets Used**

For Stage 3, we utilized datasets centered on medical instruction-following, comprehension, reasoning, and report summarization. The training corpus integrated a broad collection of **QA and understanding datasets**, including *MedQA*, *PubMedQA*, *PMC-OA*, *MedMCQA*, *PMC-InstructQA*, *MedQuAD*, *Medical-Meadow-MedQA*, *ChatDoctor-HealthCareMagic-100k*,
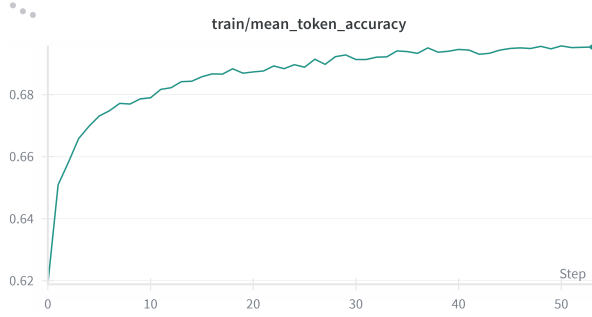
3

Figure 13. Stage 3 mean token accuracy vs. step (each step = 100 mini-batches). Accuracy increases gradually from ~0.62 to ~0.69, indicating improved instruction-following and reasoning.



Figure 14. Stage 3 training loss vs. step (each step = 100 mini-batches). Loss decreases from ~1.7 to ~1.4, showing smooth convergence and stable optimization.



Figure 15. Stage 3 output entropy vs. step (each step = 100 mini-batches). Entropy decreases from ~1.55 to ~1.38, reflecting higher model confidence and stable prediction behavior.

*AlpaCare-MedInstruct-52k*, *ChatDoctor-iCliniq*, *MedReason*, *MIMIC-IV-Ext-BHC*, *Medical-R1-Distill-Data*, *medical-o1-reasoning-SFT*, *Meadow-PubMed-Causal*, *Meadow-Medical-Flashcards*, *Meadow-MediQA*, and *Meadow-Wikidoc*. These datasets collectively provide diverse factual, reasoning, and instruction-based supervision across medical, clinical, and biomedical contexts.

In addition, we incorporated **summarization and clinical**

**reporting datasets** such as *Medical-Meadow-Cord19*, and *mimic-ext-bhc*. These datasets focus on long-form radiology and biomedical report synthesis, improving contextual understanding, summarization, and domain-specific narrative generation.

Together, this combined corpus strengthens MedMO's instruction-tuned reasoning, factual grounding, and text–image comprehension, enabling robust performance across diverse medical instruction and report-generation scenarios.

## B.4. Stage 4: Reinforcement Learning (Spatial Grounding)

### Parameters Details

- Reward functions: Label accuracy, bounding-box IoU ($\Delta$), tag count, and soft-overlong-punishment.
- Image resolution: dynamic (no fixed resize or bounding-box rescaling).
- Epsilon (policy perturbation) = 0.15.
- Epsilon_high (upper bound) = 0.25.
- Number of training epochs = 2.
- Number of batch size = 2.
- Gradient accumulation steps = 4.
- Number of generations per prompt = 8.
- Maximum prompt length = 2048 tokens.
- Maximum completion length = 1024 tokens.

### Implementation & Reproducibility Notes

- Optimizer: AdamW with default betas (0.9, 0.999) and weight decay = 0.1.
- Warm-up steps = 10% of total training steps per stage.
- Seed: All runs initialized with a fixed seed (e.g., 42) per stage; randomness only arises from data shuffling and augmentations.

### Training Dynamics

During Stage 4, MedMO was trained with reinforcement learning using the DAPO [104] algorithm to refine its spatial grounding and bounding-box localization capabilities. Each global step aggregates multiple rollouts sampled per instruction prompt. As shown in Figure 16, the bounding-box reward rises sharply from nearly zero to ~0.45 within the first 100 steps, indicating rapid adaptation of the policy to spatial localization signals. Beyond this point, the mean reward curve (blue) stabilizes around 0.42–0.45 with moderate oscillations, while the smoothed trend (red) shows a consistent upward trajectory, reflecting incremental performance gains and robust reward optimization. The steady variance band (rolling standard deviation) demonstrates that exploration remains controlled throughout training, preventing reward collapse or policy drift. Overall, the DAPO stage successfully enhances the model's spatial precision and stability in bounding-box generation tasks such as bacteria and

Figure 16. **DAPO training progress for bounding-box detection.** Mean bounding-box reward (blue) with ± rolling standard deviation (shaded) and smoothed trend (red). The consistent upward trajectory indicates effective policy optimization and stable improvement in spatial localization accuracy.
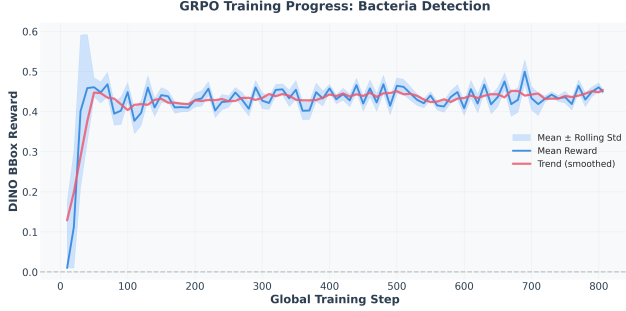
lesion detection.

### Datasets Used

For Stage 4, we utilized datasets providing explicit spatial supervision and precise bounding-box annotations for medical object detection and grounding tasks. These include *NIH Chest X-ray*, *DeepLesion*, *Bacteria Segmentation*, *CTC (Cell Tracking Challenge)*, *SLAKE*, *GRAZPEDWRI-DX*, and *MedSG*, which collectively cover anatomical structures, lesions, and microscopic cellular regions. The DAPO objective leverages bounding-box IoU and label-accuracy rewards derived from these datasets to iteratively refine spatial alignment and improve localization precision. This stage significantly enhances MedMO's visual grounding ability, leading to robust disease localization and fine-grained spatial reasoning across diverse medical modalities.

## C. Dataset Collection

We curated a unified multimodal corpus comprising **45 datasets** spanning radiology, pathology, ophthalmology, dermatology, and surgical imaging, totaling more than **26M samples**. At the core lies the **MedTrinity** dataset [96], which contributes **18.5M** publicly available instruction-following pairs. This large-scale collection integrates both image–text and text-only medical data, enabling tasks such as captioning, visual question answering (VQA), clinical reasoning, and visual grounding.

The model was trained through four progressive stages. In **Stage 1**, we used the MedTrinity dataset to establish foundational multimodal understanding across diverse imaging modalities. **Stage 2** incorporated additional VQA, grounding, and captioning datasets, and trained the model with high-resolution medical images to enhance visual reasoning and fine-grained spatial grounding.

**Stage 3** focused on medical text-only instruction data to strengthen clinical knowledge and language understand-

ing. Finally, **Stage 4** employed reinforcement learning with bounding-box supervision to further refine localization and grounding capabilities.

The datasets encompass a broad spectrum of imaging modalities (X-ray, CT, MRI, ultrasound, optical, and nuclear imaging) and biological systems (chest, brain, heart, liver, kidney, eye, colon, and tissue), ensuring comprehensive anatomical and modality coverage. For grounding supervision, we incorporated datasets containing bounding-box annotations, including *NIH Chest X-ray*, *DeepLesion*, *Bacteria*, *Wrist X-ray (boneanomaly, fracture etc.)*, *CT*, and *Cell Microscopy (DeepCell)*. This diverse corpus collectively supports robust multimodal alignment, spatial reasoning, and medical instruction tuning.

Table 8 summarizes the datasets used in MedMO's training pipeline, grouped according to their primary role in each stage.

*Note.* Several other publicly available datasets such as *TCGA* [86], *VALSET* [85], *MAMA-MIA* [28], *LLD-MMRI* [52], *CPD* [89], *CISC* [27], *CT-RATE* [29], *KIPA22* [102], and *PTCGA* [38] are already included in MedTrinity and were not trained on separately.

## D. Qualitative Results

To complement the quantitative analyses presented in the main text, Figures 17–20 provide qualitative insights into our method's performance across diverse medical imaging scenarios. These visualizations illustrate representative predictions, highlighting both successful cases and challenging examples under varied clinical conditions.

## E. Overall Training Summary

Across the four stages, MedMO progressively improves from general multimodal alignment (Stage 1) to high-resolution spatial reasoning and grounding (Stage 2), instruction-tuned language understanding (Stage 3), and reinforcement-driven grounding refinement (Stage 4). Together, these stages establish a robust, domain-aware foundation model for diverse medical imaging tasks.

Table 8. Overview of datasets used in **MedMO** training. Datasets are grouped by category, each contributing to distinct training objectives such as image captioning, multimodal and text-based instruction tuning, and spatial grounding.

| Category | Datasets | Purpose / Usage |
|---|---|---|
| **Medical Caption Data** | *MedTrinity* [96], *IU-Xray* [17], *MIMIC-CXR* [37], *CheXpert* [34], *CheXpert Plus* [12], *MEDPIX-ClinQA* [82], *ROCO* [69], *ROCO-V2* [72], *FairVLMed* [54] | Used for large-scale image–text alignment, caption-based supervision, and radiology-style report modeling across diverse imaging modalities. |
| **Medical Multimodal Instruction Data** | *VQA-Med-2019* [3], *PubMed-Vision* [15], *NIH-VQA* [74], *Quilt-LLaVA-Pretrain* [75], *MIMIC-Ext-MIMIC-CXR-VQA* [7], *VQA-RAD* [41], *PathVQA* [31], *PMC-VQA* [109], *SLAKE* [49], *CT-RATE* [11] | Facilitates multimodal instruction tuning for VQA, diagnosis, reasoning, and clinical summarization, improving image–text comprehension and task-driven responses. |
| **Medical Text Instruction Data** | *MedQA* [100], *PubMedQA* [36], *PMC-OA* [48], *MedMCQA* [67], *PMC-InstructQA* [109], *MedQuAD* [10], *Medical-Meadow-MedQA* [35], *ChatDoctor-HealthCareMagic-100k* [42], *AlpaCare-MedInstruct-52k* [108], *ChatDoctor-iCliniq* [44], *MedReason* [94], *MIMIC-IV-Ext-BHC* [1], *Medical-R1-Distill-Data* [13], *medical-o1-reasoning-SFT* [13], *Meadow-PubMed-Causal* [103], *Meadow-Medical-Flashcards* [56], *Meadow-MediQA* [70], *Meadow-Wikidoc* [57], *Medical-Meadow-Cord19* [90], *mimic-ext-bhc* [2] | Provides text-only instruction and QA supervision to enhance factual reasoning, language understanding, and medical knowledge grounding across clinical and biomedical contexts. |
| **Medical Bounding Box Data** | *NIH Chest X-ray* [26], *DeepLesion* [98], *GRAZPEDWRI-DX* [59], *SLAKE* [49], *Cell Microscopy (DeepCell, Bacteria, CTC)* [87], *MedSG* [105] | Provides explicit spatial grounding and disease-localization supervision with bounding-box annotations, enabling fine-grained object detection and improved spatial reasoning in medical imagery. |



**Ground Truth :** The heart is normal in size. The mediastinal and hilar contours are unremarkable. There is no pleural effusion or pneumothorax. The lungs appear clear. No evidence of acute process.

**MedMO :** The lungs are clear. The cardiomediastinal silhouette is within normal limits. No acute osseous abnormalities. No acute cardiopulmonary process.

Figure 17. **Qualitative result example #1.** We show model predictions compared against ground truth annotations. The input medical image is displayed on the left, with corresponding text outputs on the right. Correct predictions are highlighted (highlighted in green) to demonstrate alignment with clinical ground truth, while differences indicate areas for potential improvement.

**Ground Truth :** Frontal and lateral chest radiographs again demonstrate a vascular stent and surgical clips. Moderate to severe cardiomegaly is unchanged. There is no definite focal consolidation. Pleural and parenchymal scarring have been more fully evaluated by CTA of the chest of . A small right pleural effusion is seen. There is no appreciable pneumothorax. The visualized upper abdomen is unremarkable. Rightward deviation of the trachea is consistent with left lobe thyroid enlargement. No definite focal consolidation. Small right pleural effusion.

**MedMO :** There is a stent in the right brachiocephalic vein. There is a small right pleural effusion. There is no focal consolidation or pneumothorax. There is mild cardiomegaly. There is no pulmonary edema. There is a small hiatal hernia. There is a right upper lobe nodule, better seen on the prior CT. Small right pleural effusion. No focal consolidation.

Figure 18. **Qualitative result example #2.** We show model predictions compared against ground truth annotations. The input medical image is displayed on the left, with corresponding text outputs on the right. Correct predictions are highlighted (highlighted in green) to demonstrate alignment with clinical ground truth, while differences indicate areas for potential improvement.



**Ground Truth:** The image is a high-resolution CT scan of the thoracic region, specifically showing the lungs and the heart. The heart is centrally located, with the lungs on either side. The region of interest, located in the lower-middle right-center of the image, shows an area of the lung that is abnormal in appearance, characterized by ground-glass opacity, consolidation, or other atypical features indicative of a COVID-19 infection. This region is adjacent to normal lung tissue, which may suggest that the disease has not spread extensively and could be at an early or localized stage. The proximity to normal lung tissue could imply a potential for further spread or influence on adjacent areas, although the image alone does not confirm the directionality of this relationship.

**MedMO:** The image is a high-resolution computed tomography (HRCT) scan of the chest, showing the lungs, heart, and part of the spine. The lungs are the primary organs visible, with the heart centrally located between them and the spine visible posteriorly. The region of interest, located in the lower-middle right-center of the image, shows an area of abnormality within the lung tissue. This area exhibits a texture and density that differ from the surrounding lung parenchyma, suggesting the presence of a pathological process, likely related to COVID-19 as indicated by the caption. The abnormal region is adjacent to normal lung tissue, which may be either unaffected by the disease or could potentially be at risk of being affected due to proximity. The presence of the disease in this region could be indicative of an early stage of infection, which is consistent with the knowledge that COVID-19 can present with ground-glass opacities and consolidation in the lungs, often in the peripheral distribution of the disease.
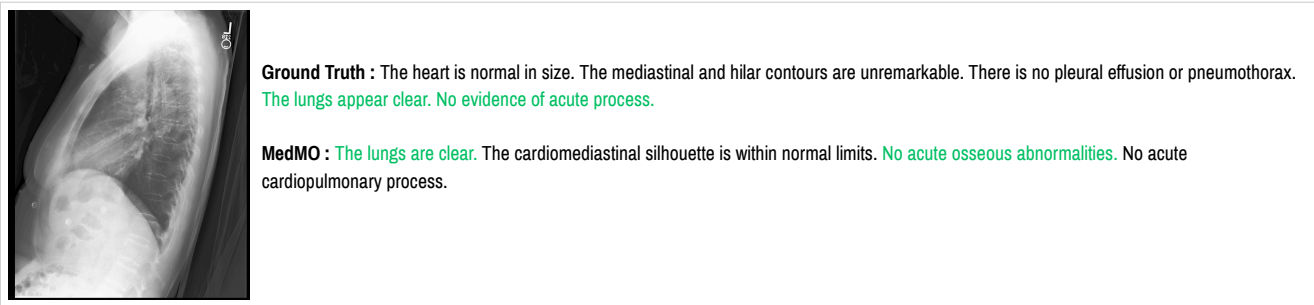
Figure 19. **Qualitative result example #3.** We show model predictions compared against ground truth annotations. The input medical image is displayed on the left, with corresponding text outputs on the right. Correct predictions are highlighted (highlighted in green) to demonstrate alignment with clinical ground truth, while differences indicate areas for potential improvement.
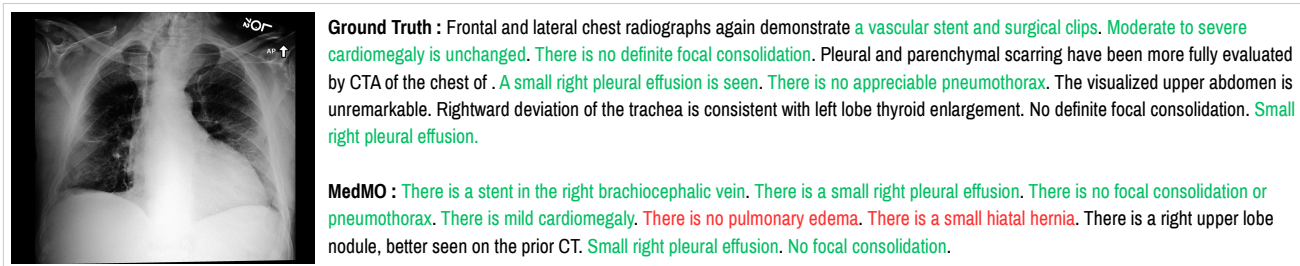


**Ground Truth:** The image is a transverse slice of a magnetic resonance imaging (MRI) scan of the brain, showing various brain structures such as the cerebral cortex, ventricles, and other intracranial structures. A region of interest is located at the center of the brain, which is the cerebral hemispheres, and it appears to be a transverse slice through the axial plane, likely at the level of the cerebral cortex, possibly involving the cerebral cortex and adjacent brain tissue. The region of interest, located centrally in the image, there is an area that appears to be abnormal, with an area ratio of 0.4% of the image, indicating a potential abnormality. This region is characterized by a difference in tissue character, which may suggest a variation in tissue density or composition, such as altered signal intensity, possibly due to the presence of a brain tumor. The abnormality is situated in the cerebral cortex, which is the central part of the brain, which could be affecting or be affected by the disease process, potentially disrupting the normal brain tissue surrounding it. The abnormality's relationship to other brain structures could be indicative of a pathological process that may be influencing or being influenced by adjacent brain tissue, such as causing a mass effect or being influenced by surrounding brain structures, potentially altering their function or being affected by them.

**MedMO:** The image is a magnetic resonance imaging (MRI) scan of the brain, showing various brain structures such as the cerebral cortex, ventricles, and possibly the basal ganglia or thalamus, without any medical devices visible. A region of interest is located at center part of the image horizontally and upper-middle part of image vertically, occupying approximately 0.4% of the total area. The region of interest, located in the cerebral cortex, exhibits an abnormality that differs in appearance from the surrounding brain tissue, potentially indicating a brain tumor. This abnormality could be affecting or be affected by adjacent brain structures due to its proximity, which may suggest an interaction or influence on nearby tissue, although the exact relationship depends on the nature of the tumor and its growth pattern. The MRI scan, with its high soft-tissue contrast, is instrumental in identifying such abnormalities.

Figure 20. **Qualitative result example #4.** We show model predictions compared against ground truth annotations. The input medical image is displayed on the left, with corresponding text outputs on the right. Correct predictions are highlighted (highlighted in green) to demonstrate alignment with clinical ground truth, whi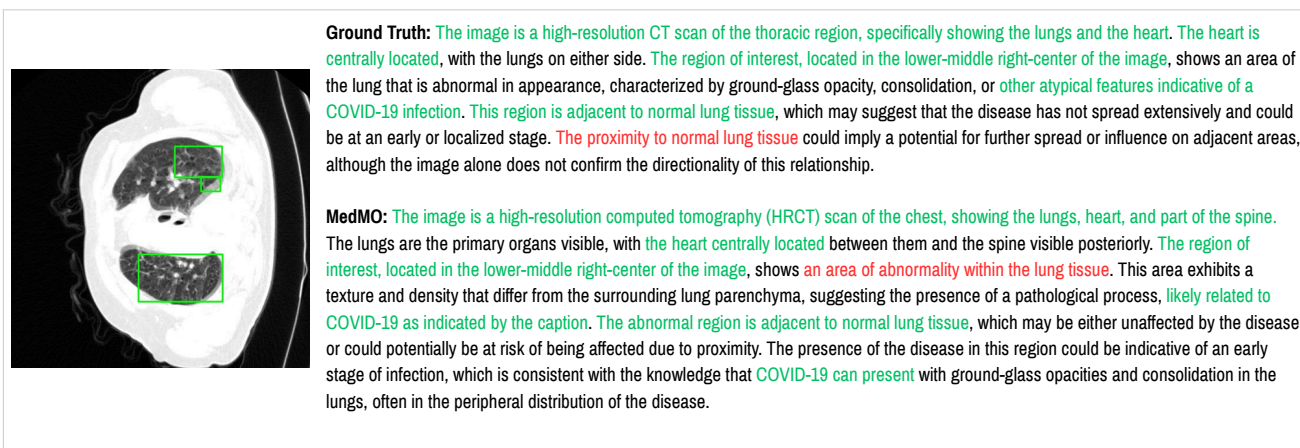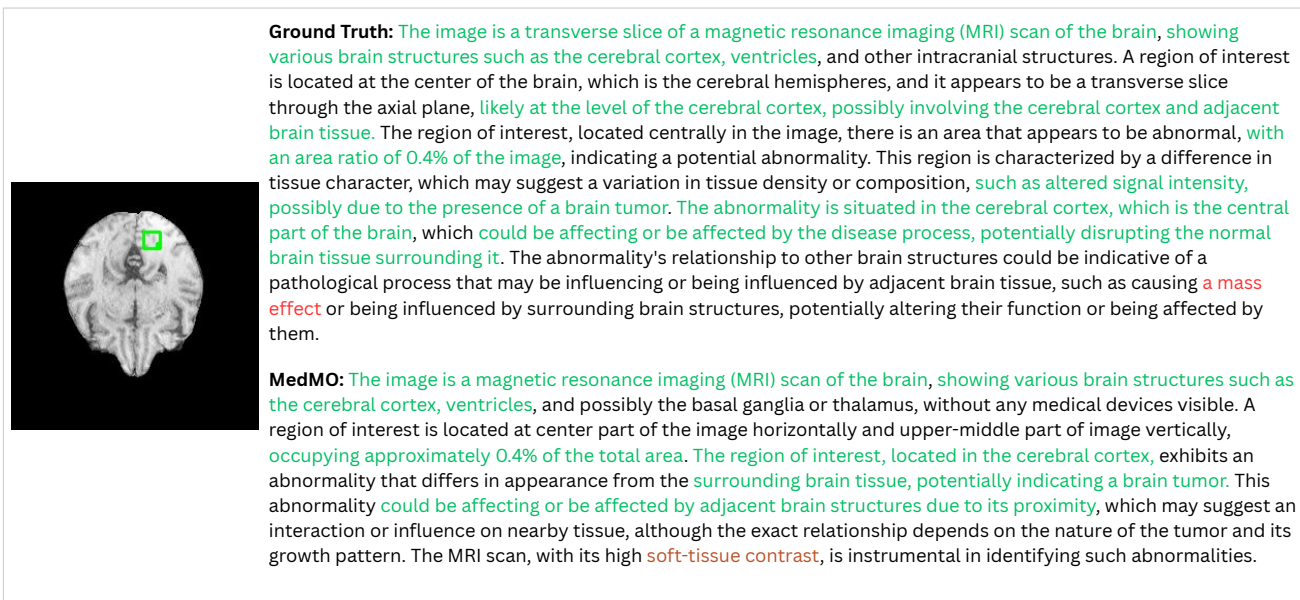le differences indicate areas for potential improvement.