

MEDIC-AD: Towards Medical Vision-Language Model’s Clinical Intelligence

Anonymous CVPR submission

Paper ID 22318

Abstract

001 *Lesion detection, symptom tracking, and visual explainability*
 002 *are central to real-world medical image analysis, yet*
 003 *current medical Vision-Language Models (VLMs) still lack*
 004 *mechanisms that translate their broad knowledge into clinically*
 005 *actionable outputs. To bridge this gap, we present*
 006 *MEDIC-AD, a clinically oriented VLM that strengthens*
 007 *these three capabilities through a stage-wise framework.*
 008 *First, learnable anomaly-aware tokens (<Ano>) encourage*
 009 *the model to focus on abnormal regions and build more dis-*
 010 *criminative lesion centered representations. Second, inter-*
 011 *image difference tokens (<Diff>) explicitly encode tempo-*
 012 *ral changes between studies, allowing the model to dis-*
 013 *tinguish worsening, improvement, and stability in disease*
 014 *burden. Finally, a dedicated explainability stage trains the*
 015 *model to generate heatmaps that highlight lesion-related*
 016 *regions, offering clear visual evidence that is consistent*
 017 *with the model’s reasoning. Through our staged design,*
 018 *MEDIC-AD steadily boosts performance across anomaly*
 019 *detection, symptom tracking, and anomaly segmentation,*
 020 *achieving state-of-the-art results compared with both closed*
 021 *source and medical-specialized baselines. Evaluations on*
 022 *real longitudinal clinical data collected from real hospital*
 023 *workflows further show that MEDIC-AD delivers stable*
 024 *predictions and clinically faithful explanations in practical*
 025 *patient-monitoring and decision-support workflows.*

026 1. Introduction

027 Vision-Language Models (VLMs) have rapidly
 028 evolved [14, 39, 61] from simple image-text tasks
 029 such as visual question answering (VQA) [2] and caption-
 030 ing [49] to more advanced capabilities including visual
 031 grounding [44, 57] and multi-image reasoning [3, 11, 33].
 032 These advances have inspired the emergence of *Medical*
 033 *Foundation VLMs* [34, 40, 47], which aim to integrate
 034 visual and textual medical knowledge for comprehensive
 035 diagnostic understanding. Trained on large-scale image-
 036 report pairs and multimodal instructions, these models have
 037 achieved strong results across tasks such as disease classi-

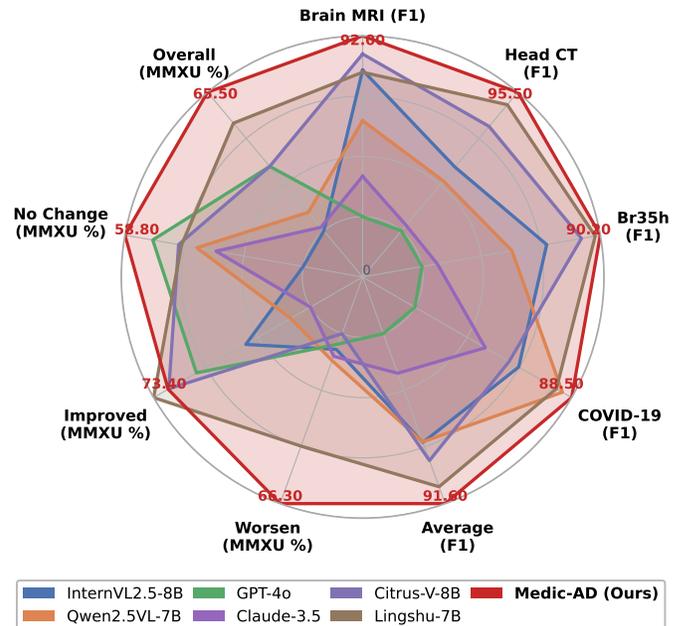


Figure 1. Overall performance of VLMs on Medical Anomaly Detection and Medical Symptom Tracking (MMXU [41]).

038 fication, report generation, and Med-VQA, demonstrating
 039 the promise of language-driven clinical reasoning.

040 However, most Medical Foundation VLMs remain op-
 041 timized for broad medical knowledge coverage rather than
 042 real clinical application [50, 56]. Their training typically re-
 043 lies on long-form captioning, OCR-based instruction tun-
 044 ing, and medical chain-of-thought reasoning that enhance
 045 generic reasoning ability, but overlook key properties re-
 046 quired for real-world clinical workflows [1]: (1) accu-
 047 rate lesion detection, (2) reliable temporal symptom track-
 048 ing, and (3) transparent visual explainability of the rea-
 049 soning process. Addressing these limitations demands a
 050 paradigm shift from generalized intelligence toward clini-
 051 cally grounded perception, and understanding.

052 To that end, we explore three research questions guiding
 053 the design of a clinically usable medical VLM.

RQ1: How can lesion and symptom recognition be improved in VLMs for real clinical settings? Even as medical VLMs expand in knowledge, accurate abnormality detection remains essential for safe deployment. We define an *abnormality* as any pathological deviation within an image and propose to enhance this recognition through explicit anomaly-aware representations. By injecting learnable $\langle \text{Ano} \rangle$ tokens into the transformer layers, the model highlights abnormal regions and strengthens its discriminative reasoning. Experiments on brain MRI, head CT, and chest X-ray datasets show that this design achieves strong performance in medical anomaly detection.

RQ2: How can a VLM disentangle temporal medical images to enable more accurate symptom tracking? Existing foundation models that support multi-image inputs typically concatenate visual features, thus failing to capture the temporal progression between scans. To model clinically meaningful changes, we introduce $\langle \text{Diff} \rangle$ tokens that compare anomaly features extracted from multiple images of the same patient. These representations allow the model to reason about whether a condition has worsened, improved, or remained unchanged. On benchmarks such as MMXU [41], which assess longitudinal symptom understanding, our approach achieves superior performance, highlighting its effectiveness as a practical clinical tool for patient monitoring.

RQ3: How can visual explainability be integrated into medical VLM reasoning? Interpretability is indispensable for clinical decision-making. To visually justify model’s predictions, we design a heatmap decoder that fuses anomaly features and visual features to generate attention maps highlighting regions responsible for each prediction. These region-level explanations enhance transparency by providing visual evidence for both lesion detection and change assessment, ultimately bridging the gap between black-box reasoning and clinical trust.

These three research directions culminate in **MEDIC-AD**, a stage-wise medical VLM with clinical intelligence, designed to integrate anomaly detection, temporal reasoning, and visual explainability. MEDIC-AD is trained in three stages: *Stage 1: Anomaly Detection*. Learn discriminative abnormality embeddings via injected anomaly tokens, $\langle \text{Ano} \rangle$, adapting contrastive architecture between normal and abnormal regions for enhancing sensitivity to pathological cues. *Stage 2: Difference Reasoning*. Encode cross-scan variations using $\langle \text{Diff} \rangle$ tokens that disentangle temporal progression of abnormal features and enable fine-grained symptom tracking. *Stage 3: Visual Explainability*. Generate interpretable heatmaps that ground textual outputs on visually abnormal regions, ensuring verifiable reasoning.

Through extensive evaluation, MEDIC-AD consistently demonstrates state-of-the-art (SOTA) performance across diverse medical modalities and tasks as shown in Fig. 1. It outperforms medical foundation models [34, 50, 56], anomaly-specialized models [17, 55], and closed-source counterparts [4, 24] in both lesion detection, and temporal symptom reasoning. Moreover, MEDIC-AD delivers superior visual interpretability generating spatially grounded explanations that align model decisions with clinical evidence. Beyond numerical gains, its stage-wise design encodes the clinical diagnostic workflow—detect, compare, explain—into the model’s learning curriculum, transforming general-purpose vision-language understanding into clinically actionable intelligence.

Our main contributions are as follows:

- We present a unified, stage-wise framework (MEDIC-AD) that integrates anomaly detection, longitudinal reasoning, and visual grounding to enable explainable medical inference.
- We introduce anomaly- and difference-token mechanisms that endow medical VLMs with explicit lesion sensitivity and temporal reasoning capability.
- We conduct comprehensive evaluations on multiple medical tasks, as well as on real longitudinal datasets from hospital sites, demonstrating superior reliability and usability compared to both open- and closed-source foundation models, and showcasing deployment readiness for real-world clinical practice workflows.

2. Related Works

2.1. Vision-Language Models for Medical Imaging

Vision Language Models (VLMs) [3, 11, 14, 33, 39, 61] have unified visual and textual reasoning across domains through large-scale contrastive or instruction-tuned learning. Building upon these, medical VLMs have adapted multimodal alignment to clinical imaging and reporting tasks. Early medical VLMs focused on contrastive alignment and report-level representation learning [6, 51], while later instruction-tuned architectures expanded multimodal reasoning through large scale medical-text pretraining [34, 46, 48, 54]. More recently, Lingshu and Citrus-V [50, 56], built on Qwen-VL 2.5 [3], introduced multi-stage training with shallow/deep alignment, medical instruction tuning, and reinforcement learning with verifiable rewards, achieving state-of-the-art results on single-image medical VQA benchmarks such as SLAKE, PathVQA, VQA-RAD, and OmniMedVQA [19, 22, 32, 48]. These datasets collectively evaluate anatomical localization, factual consistency, and report generation, forming the empirical foundation for single-image medical VLMs. Beyond single-image understanding, medical VLMs have advanced toward difference-aware reasoning, modeling longitudinal changes and dis-

156 ease progression between paired studies. Generic differ-
 157 ence captioning frameworks describe visual changes across
 158 image pairs [43, 58], whereas medical-specific longitudi-
 159 nal reasoning integrates anatomical or report-based tempo-
 160 ral modeling [12, 21]. A recent unified framework [15]
 161 proposes a *Report Generator–Answer Generator (RG–AG)*
 162 architecture. While these studies have achieved meaning-
 163 ful progress in understanding longitudinal medical images,
 164 they remain largely task-specific, focusing on objectives
 165 such as VQA and report generation, without fully leverag-
 166 ing the extensive medical knowledge and generative reason-
 167 ing capabilities offered by Medical Foundation VLMs.

168 2.2. Zero-Shot Anomaly Detection

169 Traditional anomaly detection (AD) methods primarily fo-
 170 cused on low-level visual irregularities in industrial datasets
 171 such as MVTec-AD [8] and VisA [62], later extending to-
 172 ward zero-shot recognition through text–image alignment.
 173 CLIP-based approaches introduced adapter- or prompt-
 174 based mechanisms for open-vocabulary detection [10, 23,
 175 60], while Q-Former-based [36] architectures further con-
 176 nected anomaly detection and instruction tuning [17, 55]
 177 in general AD tasks. In the medical context, unified
 178 benchmarks such as BMAD [7] integrate diverse medical
 179 anomaly detection datasets—covering Brain MRI, Chest X-
 180 Ray, Liver CT, Retinal OCT, and Pathology—into a sin-
 181 gle evaluation framework. BMAD consolidates various
 182 modality-specific datasets to enable consistent cross-dataset
 183 and cross-organ evaluation under a unified protocol. In par-
 184 allel, chest-specific datasets such as ChestX-Det [38] fur-
 185 ther provide detailed pixel-level annotations for thoracic
 186 disease localization and anomaly segmentation. Collec-
 187 tively, these works mark a shift from handcrafted features
 188 to medically grounded anomaly understanding.

189 2.3. Explainability in Vision–Language Models

190 Explainability has become a key criterion for assessing the
 191 reliability of VLMs, especially in safety-critical domains
 192 such as medicine. Recent VLMs utilize cross-attention
 193 heatmaps and token-conditioned activations to visualize
 194 how linguistic tokens attend to visual regions during reason-
 195 ing, thereby revealing the internal correspondence between
 196 textual semantics and spatial evidence [16, 35, 39, 42]. Such
 197 mechanisms improve transparency and enable systematic
 198 model auditing and error analysis by linking visual atten-
 199 tion patterns with generated textual outputs.

200 In medical imaging, explainability has been advanced
 201 through explicit grounding and concept-level alignment.
 202 Grounded VLMs explicitly associate textual rationales
 203 with anatomical regions [6, 15, 25], while concept-
 204 disentanglement approaches align clinical entities with vi-
 205 sual concepts to enhance interpretability and trustworthi-
 206 ness [37]. Building upon these developments, we extend

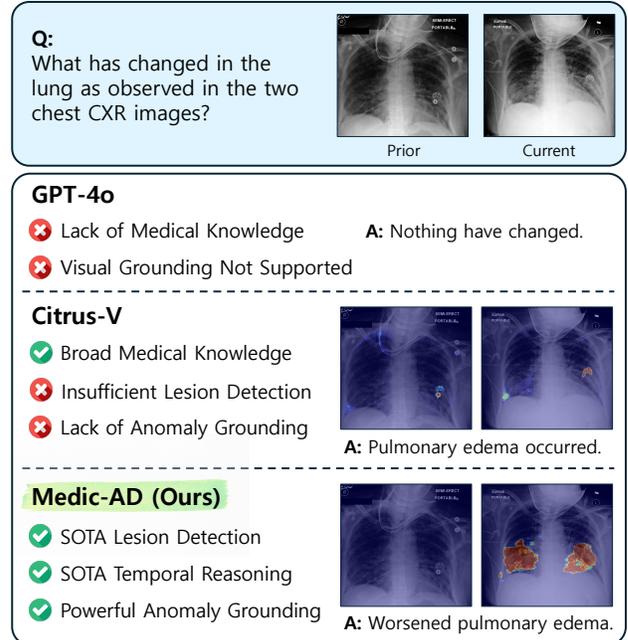


Figure 2. Comparison of VLMs on clinical applications. Medic-AD provides stronger lesion detection, temporal reasoning, and visual grounding than GPT-4o [24] and Citrus-V [50].

visual grounding beyond a mere auxiliary visualization tool. In our medical VLM, interpretability is achieved by grounding the reasoning process on anomalous features such as lesions or symptoms, thereby providing explicit visual evidence that supports the model’s clinical conclusions and ensuring clinically verifiable visual explainability.

3. Methodology

3.1. Overview

Standard VLMs encode an input image \mathbf{I} and textual instruction \mathbf{T} into a joint multimodal sequence to generate a response \mathbf{R} :

$$\mathbf{R} = f_l([f_p(\mathbf{V}); \text{Emb}(\mathbf{T})]), \quad (1)$$

where $f_p(\cdot)$ denotes a visual projection layer that maps visual features \mathbf{V} , extracted from \mathbf{I} via a vision encoder, into the text embedding space, and $f_l(\cdot)$ represents the large language model (LLM). While this general formulation supports broad multimodal reasoning, it lacks the inductive biases required for clinically meaningful perception including reliable lesion localization, temporal tracking, and visual justification, all of which are essential for trustworthy decision support.

MEDIC-AD extends this paradigm into a clinically grounded reasoning framework through a stage-wise optimization pipeline, composed of three progressive stages where each stage incrementally enhances the model’s reasoning capability. Stage 1 learns anomaly-aware represen-

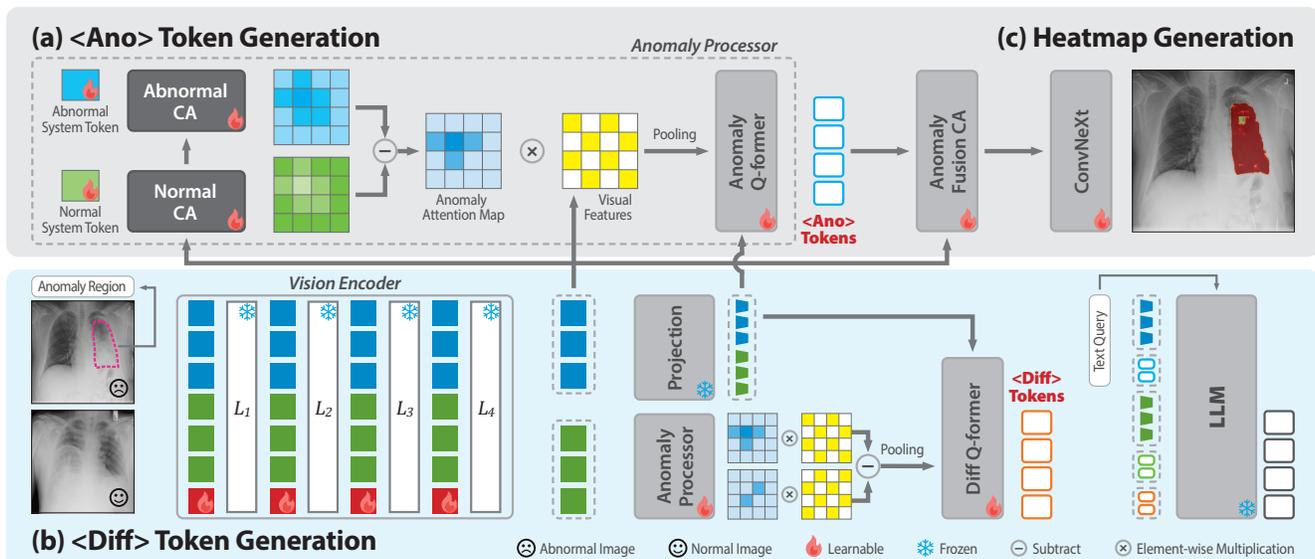


Figure 3. Architecture of MEDIC-AD. (a) Stage 1: <Ano> Token Generation, (b) Stage 2: <Diff> Token Generation, and (c) Stage 3: Heatmap Generation illustrate each stage of the proposed framework. Note that CA denotes Cross-Attention.

233 tations that encode lesion-specific semantics. Stage 2 builds
 234 on these representations to disentangle temporal variations
 235 between prior and current studies, yielding difference tokens.
 236 Stage 3 introduces grounding supervision that aligns
 237 the learned anomaly features with spatial heatmaps, en-
 238 abling visually verifiable predictions.

239 We first introduce anomaly-aware tokens, <Ano>, de-
 240 rived from a cross-attention mechanism applied to the vi-
 241 sually enhanced feature representation \mathbf{V}^* . Here, \mathbf{V}^*
 242 denotes the anomaly-augmented visual features obtained by
 243 incorporating *visual soft prompts* [26] into the original vi-
 244 sual embeddings \mathbf{V} . This process encourages the model to
 245 emphasize lesion-relevant regions and discriminative cues.
 246 The resulting tokens are concatenated with the visual and
 247 textual embeddings as

$$248 \quad \mathbf{R} = f_l([f_p(\mathbf{V}^*); \langle \text{Ano} \rangle; \text{Emb}(\mathbf{T})]). \quad (2)$$

249 Next, to model temporal changes between prior and cur-
 250 rent images, our model learns anomaly-aware representa-
 251 tions across time and produces <Diff> tokens. Given two
 252 input images, their corresponding anomaly-augmented vi-
 253 sual features, \mathbf{V}_1^* and \mathbf{V}_2^* , are extracted by the vision
 254 encoder and formulated as

$$255 \quad \mathbf{R} = f_l([f_p(\mathbf{V}_1^*); \langle \text{Ano} \rangle; f_p(\mathbf{V}_2^*); \langle \text{Ano} \rangle; \\ \text{Emb}(\mathbf{T}); \langle \text{Diff} \rangle]), \quad (3)$$

256 following the modified chat template as illustrated in Ap-
 257 pendix Sec. A. Finally, the anomaly-aware tokens, <Ano>,
 258 are fed into a heatmap decoder $f_h(\cdot)$ together with the cor-
 259 responding visual features \mathbf{V}^* to generate grounding maps
 260 \mathbf{M} . These maps provide region-level visual evidence that

supports and justifies textual predictions as

$$\mathbf{M} = f_h([f_p(\mathbf{V}^*); \langle \text{Ano} \rangle]). \quad (4)$$

3.2. Architecture and Stage-wise Training

In this section, we present the detailed architecture of MEDIC-AD and the corresponding training pipelines designed to implement the framework illustrated in Sec. 3.1. Each stage progressively enhances the capability of the baseline Medical Foundation VLM, Lingshu [56], which is a strong backbone model pretrained on large-scale medical data, by introducing specialized modules for anomaly reasoning, temporal difference analysis, and visual explainability.

Stage 1: Anomaly Detection. The first stage focuses on training an *Anomaly Processor* that produces <Ano> tokens, compact latent representations capturing lesion-related semantics, as illustrated in Fig. 3 (a). These tokens are constructed through two learnable system tokens, the *Abnormal System Token* and *Normal System Token*. They interact with multi-scale visual features extracted from four intermediate layers of the vision encoder via a cross-attention mechanism, producing *Abnormal* and *Normal Attention Scores* for each visual patch.

To preserve the pretrained vision encoder’s representational stability while adapting it for anomaly detection, we adopt *Visual Soft Prompt Tuning* [26] to the selected four layers instead of fully updating their parameters. Unlike conventional attention head using Softmax normalization, our design applies Sigmoid activation to obtain patch-wise anomaly probabilities. The difference between abnormal

290 and normal attention weights yields an *Anomaly Attention*
291 *Map*, which reflects the likelihood of each patch being ab-
292 normal. This map modulates the original visual features
293 through element-wise multiplication, adjusting their mag-
294 nitudes according to anomaly salience.

295 Subsequently, 2D global pooling is performed over
296 the modulated visual features to derive *Anomaly Queries*.
297 These queries are passed through an *Anomaly Q-Former*,
298 where the LLM-projected visual tokens act as keys and
299 values. The Anomaly Q-Former outputs are then fed into
300 a 2-layer MLP to yield the final Anomaly-aware tokens,
301 $\langle \text{Ano} \rangle$, which are used in both downstream LLM inference
302 and later heatmap generation in Stage 3.

303 Training for this stage utilizes a diverse collection of
304 medical anomaly datasets spanning MRI, X-ray, and CT
305 modalities, including **BMAD**, **ChestX-Det** [7, 38], as well
306 as multimodal VQA datasets such as **SLAKE**, **PathVQA**,
307 and **VQA-RAD** [19, 32, 48], ensuring both robust visual
308 grounding and generalizable medical reasoning capability.

309 **Stage 2: Difference Reasoning.** The second stage fo-
310 cuses on modeling inter-image differences to analyze dis-
311 ease progression over time. Here, the goal is to learn dif-
312 ference tokens, $\langle \text{Diff} \rangle$, that *disentangle* the variations in
313 abnormal regions across time or paired studies (e.g., follow-
314 up vs. baseline scans), effectively separating genuine patho-
315 logical progression from visual or acquisition-related noise.

316 As shown in Fig. 3 (b), the modulated visual features
317 derived in Stage 1 from two images are contrasted and dis-
318 entangled through a *Diff Q-Former*, which isolates lesion-
319 specific change patterns. Then, each image’s projected vi-
320 sual tokens, $f_p(\mathbf{V}_1^*)$ and $f_p(\mathbf{V}_2^*)$, serves as keys and val-
321 ues to encode structured inter-image relationships. Passing
322 these through the Diff Q-Former and a subsequent 2-layer
323 MLP yields the difference tokens, $\langle \text{Diff} \rangle$, which are ap-
324 pended to the multimodal input sequence. By explicitly iso-
325 lating temporal anomalies from static visual context, this
326 stage enables the LLM to reason over fine-grained temporal
327 variations in lesion appearance or intensity, thereby enhanc-
328 ing its ability for longitudinal disease reasoning.

329 Stage 2 training requires temporally paired or longitu-
330 dinal datasets. We use **MIMIC-Diff-VQA**[21], a dataset
331 built for multi-image reasoning in clinical follow-up sce-
332 narios, allowing the model to learn spatial correspondence
333 and temporal progression patterns in real patient studies.

334 **Stage 3: Visual Explainability.** The final stage intro-
335 duces a heatmap generation module designed to achieve vi-
336 sual grounding and enhance interpretability. While prior
337 Medical Foundation VLMs often rely on pretrained vision
338 decoders (e.g., SAM2 [45] used in Citrus-V [50]), our ap-
339 proach leverages the learned $\langle \text{Ano} \rangle$ tokens and ConvNeXt-

based [52] segmentation head to directly link visual reason-
ing with interpretable evidence.

340 As illustrated in Fig. 3 (c), we fuse $\langle \text{Ano} \rangle$ tokens with
341 the vision encoder’s intermediate feature maps via a fusion
342 block, reinforcing the model’s focus on lesion-relevant re-
343 gions driving the LLM’s prediction. The fused features are
344 then processed by a compact ConvNeXt-based segmenta-
345 tion head to generate a heatmap \mathbf{M} spatially aligned with
346 the input image. This heatmap is overlaid on the original
347 image to provide region-level visual evidence that supports
348 the textual output, thereby connecting model reasoning with
349 clinically observable cues.

350 Stage 3 is trained on datasets with pixel-level segmen-
351 tation masks, such as selected subsets of **BMAD** and
352 **ChestX-Det** [7, 38]. Leveraging anomaly-token-guided fu-
353 sion, our model achieves substantially improved anomaly-
354 localization accuracy compared with recent grounding-
355 based medical VLMs, as demonstrated in Sec. 4.3. For the
356 more detailed training configurations of each stages, please
357 see Appendix Sec. B.

360 4. Experiments

361 To validate the effectiveness of MEDIC-AD, we con-
362 duct a series of experiments corresponding to each stage
363 of the proposed framework introduced in Sec. 3. Each
364 stage is validated on a task specifically aligned with its
365 objective. Stage 1 evaluates the discriminative capabil-
366 ity of the learned $\langle \text{Ano} \rangle$ tokens via **Medical Zero-shot**
367 **Anomaly Detection**. Stage 2 assesses temporal reason-
368 ing through the **MMXU Benchmark** for medical symp-
369 tom tracking. Finally, Stage 3 examines **Medical Visual**
370 **Explainability**, evaluating region-level grounding and the
371 consistency between visual evidence and textual reasoning
372 using segmentation-based metrics.

373 As this section focuses on the stage-specific tasks cen-
374 tral to our framework, evaluations on conventional med-
375 ical VQA benchmarks (**VQA RAD** [32], **SLAKE** [48],
376 **PathVQA** [19], and **MMMU Med** [59]) are provided in
377 Appendix Sec. C, demonstrating that MEDIC-AD success-
378 fully preserves general medical knowledge of backbone
379 while simultaneously reinforcing its stage-wise clinical ap-
380 plicability.

381 4.1. Medical Zero-shot Anomaly Detection

382 **Experimental Settings.** In the zero-shot anomaly detec-
383 tion setting, the model is tested on datasets that are entirely
384 unseen during training to ensure a fair comparison with
385 competing models. We evaluate across four heterogeneous
386 modalities—**Brain MRI**, **Head CT**, **Br35h**, and **COVID-**
387 **19 X-ray** [13, 18, 28, 30]—to assess generalization capabil-
388 ity. Each test sample is queried with a consistent instruction
389 prompt: “*Is there any abnormality in this image?*” The

Table 1. Results on Zero-shot Medical Anomaly Detection (Brain MRI [28], Head CT [30], Br35h [18], and COVID-19 [13]). For each dataset, we report Precision, Recall, and F1. The rightmost column shows the average F1 over all tasks.

Category	Model	Size	Brain MRI			Head CT			Br35h			COVID-19			Avg. F1
			Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1	
General	Qwen2.5-VL	7B	76.1	92.9	83.6	61.6	100.0	76.3	66.4	94.5	78.0	88.6	83.4	85.9	81.0
	InternVL2.5	8B	81.2	97.4	88.6	67.6	96.0	79.3	71.2	99.0	82.8	94.8	60.2	73.6	81.1
Closed	GPT-4o	–	59.3	98.6	74.1	48.7	100.0	65.5	48.9	99.9	65.6	29.1	92.9	44.4	62.4
	Claude-3.5	–	64.0	100.0	78.1	50.0	100.0	66.7	51.2	100.0	67.7	47.2	100.0	64.2	69.2
Anomaly	AnomalyGPT	13B	61.0	96.8	74.8	50.8	97.0	66.7	47.2	88.6	61.6	33.9	57.2	42.5	61.4
	Anomaly-OV	7B	53.1	71.6	61.0	39.8	66.0	49.6	42.2	73.0	53.5	29.6	47.0	36.3	50.1
Medical	LLaVA-MED	7B	69.0	95.5	80.1	54.1	77.9	63.8	60.8	92.7	73.4	46.1	86.3	60.1	69.4
	Citrus-V	8B	85.5	95.5	90.2	78.7	100.0	88.1	79.5	97.5	87.6	58.4	90.4	70.9	84.2
	Lingshu	7B	83.1	94.3	88.4	91.4	94.1	92.8	86.0	93.3	89.5	84.0	84.4	84.2	88.7
	MEDIC-AD	7B	91.1	92.9	92.0	95.1	96.0	95.5	90.8	89.5	90.2	96.6	81.6	88.5	91.6

Table 2. Results on MMXU [41]. Models are categorized into general-purpose, closed-source, and medical-domain VLMs.

Category	Model	Size	Worsen	Improved	No Change	Overall (↑)
General	InternVL2.5	8B	0.486	0.607	0.402	0.498
	Qwen2.5-VL	7B	0.499	0.545	0.513	0.519
Closed	Claude-3.5	–	0.494	0.518	0.493	0.502
	GPT-4o	–	0.480	0.675	0.559	0.571
Medical	Citrus-V	8B	0.468	0.713	0.532	0.571
	Lingshu	7B	0.597	0.734	0.529	0.620
	MEDIC-AD (Ours)	7B	0.663	0.714	0.588	0.655

390 model’s binary response is evaluated using the F1 score to
391 capture both precision and recall performance.

392 We compare MEDIC-AD against a range of baselines
393 spanning three categories: (1) **General-purpose open-**
394 **source VLMs:** Qwen2.5-VL-7B and InternVL2.5-8B [3,
395 11], (2) **Closed-source models:** GPT-4o and Claude-3.5 [4,
396 24], and (3) **Anomaly Detection Specialized VLMs:**
397 AnomalyGPT and Anomaly-OV [17, 55], as well as (4)
398 **Medical Foundation VLMs:** LLaVA-Med-7B, Citrus-V-
399 8B, and Lingshu-7B [34, 50, 56].

400 **Results.** As summarized in Table 1, MEDIC-AD achieves
401 **SOTA performance** across all four datasets, demonstrating
402 superior generalization to unseen medical imaging modal-
403 ities and conditions. The results indicate that the learned
404 anomaly-aware tokens, $\langle A_{no} \rangle$, effectively capture lesion-
405 relevant features, enabling reliable and robust zero-shot
406 abnormality discrimination even in unseen datasets. No-
407 tably, MEDIC-AD surpasses all medical-specialized base-
408 lines and even outperforms closed-source models, vali-
409 dating the efficacy of its anomaly representation learning.
410 While some closed-source models (e.g., GPT-4o) occasion-
411 ally produce conservative outputs for normal cases due to
412 safety-alignment bias [29, 53] (e.g., responding with “I’m
413 unable to analyze medical images ...”), MEDIC-AD consis-
414 tently produces well-calibrated predictions across both nor-

mal and abnormal cases, indicating greater clinical reliabil- 415
ity and decision consistency. Overall, Stage 1 demonstrates 416
that anomaly representation learning offers a reliable founda- 417
tion for zero-shot medical abnormality detection. 418

4.2. Medical Symptom Tracking 419

Experimental Settings. We evaluate temporal reasoning 420
ability using the MMXU benchmark [41], which involves 421
paired chest X-ray studies from the same patient. For each 422
instance, the model must classify disease progression as 423
worsened, *improved*, or *unchanged* based on two images 424
and a multiple-choice question such as: “What is the condi- 425
tion of the left lower lung zone across the two chest CXR 426
images? A: No significant change, B: Improved, C: Wors- 427
ened.” Only models supporting multi-image reasoning are 428
included in this comparison, and the evaluation metric fol- 429
lows the accuracy of categorical predictions derived from 430
the model’s textual responses. 431

Results. As shown in Tab. 2, MEDIC-AD demonstrates 432
clear advantages over existing multimodal models on the 433
MMXU benchmark. By leveraging the $\langle Diff \rangle$ tokens 434
representations disentangled through the Diff Q-Former, the 435
model effectively captures localized lesion changes while 436
being robust to irrelevant appearance variations such as il- 437
lumination or positional shifts. Unlike general VLMs that 438

Table 3. Visual grounding performance on BMAD [7] (BraTS2021 [5], RESC [20], BTCV + LiTs [9, 31]) and ChestX-Det [38] datasets. Each dataset reports AUC and mIoU metrics (higher is better).

Model	BraTS2021		RESC		BTCV + LiTs		ChestX-Det	
	AUC (\uparrow)	mIoU (\uparrow)	AUC (\uparrow)	mIoU (\uparrow)	AUC (\uparrow)	mIoU (\uparrow)	AUC (\uparrow)	mIoU (\uparrow)
Citrus-V	98.8	32.6	87.6	2.1	82.1	1.7	98.0	12.4
Medic-AD (Ours)	99.8	87.6	100	97.2	97.2	83.6	99.8	79.8

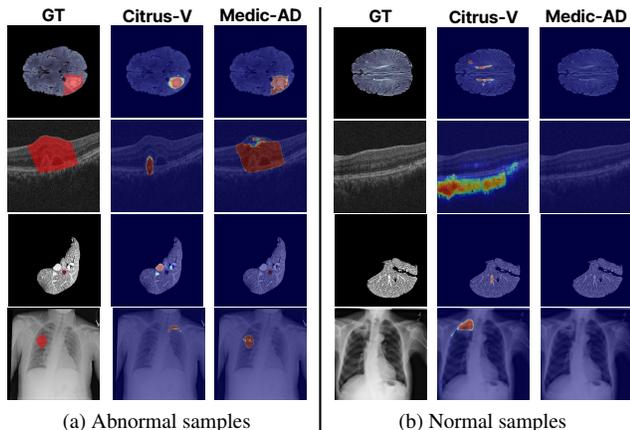


Figure 4. Visual Grounding comparison between MEDIC-AD and Citrus-V [50] on diverse abnormal and normal samples.

439 simply concatenate multiple images, MEDIC-AD explicitly
 440 encodes inter-image relationships, yielding consistent
 441 reasoning about temporal dynamics. In particular, qualitative
 442 inspection shows that MEDIC-AD highlights pathological
 443 regions with true clinical changes (e.g., consolidation
 444 growth or opacity reduction), whereas other models
 445 often mistake global contrast shifts for disease progression
 446 (Fig.2). These findings indicate that Stage 2 effectively
 447 separates clinically relevant temporal changes.

448 4.3. Medical Visual Explainability

449 **Experimental Settings.** To assess visual grounding and
 450 interpretability, we evaluate MEDIC-AD on medical
 451 datasets that include pixel-level anomaly masks, specifically
 452 a subset of **BMAD** [7] (BraTS2021 [5], RESC [20],
 453 and BTCV + LiTs [9, 31]) and the **ChestX-Det** [38] dataset.
 454 For each image, the model generates a heatmap using the
 455 same anomaly-detection query as in Sec. 4.1. Predicted
 456 heatmaps are compared with ground-truth masks using
 457 AUC and mIoU. We compare against **Citrus-V** [50], the
 458 most recent medical VLMs supporting visual grounding.

459 **Results.** MEDIC-AD shows consistently strong performance
 460 across datasets, outperforming Citrus-V on both
 461 AUC and mIoU (Tab. 3). By integrating $\langle \text{Ano} \rangle$ tokens
 462 with intermediate visual features, MEDIC-AD produces

463 heatmaps that more accurately localize pathological regions
 464 aligned with the model’s textual rationale. In contrast,
 465 Citrus-V, which use a SAM2 decoder [45], tends to produce
 466 less precise masks, sometimes highlighting non-lesion
 467 areas or yielding diffuse activations in normal images. Representative
 468 examples in Fig. 4 highlight these differences. These results
 469 show that Stage 3 improves the alignment between visual
 470 evidence and model reasoning, yielding more interpretable
 471 and clinically coherent responses.

472 5. Analysis

473 In this section, we present a comprehensive analysis of
 474 the proposed MEDIC-AD framework through ablation studies,
 475 hyperparameter sensitivity experiments, and real-world
 476 evaluations. Across all studies, the results consistently support
 477 the central claim of this work: temporal reasoning in
 478 medical image pairs fundamentally benefits from coherent
 479 integration of anomaly-aware spatial cues and temporally
 480 grounded difference representations.

481 5.1. Effect of $\langle \text{Ano} \rangle$ Tokens on Temporal Reasoning

482 The first analysis focuses on how $\langle \text{Ano} \rangle$ tokens contribute
 483 to constructing reliable temporal representations. In MEDIC-AD,
 484 the construction of $\langle \text{Diff} \rangle$ tokens is grounded in the
 485 anomaly-aware visual features generated during the $\langle \text{Ano} \rangle$
 486 tokens estimation process, where patch-wise anomaly
 487 likelihood modulates the magnitude of visual representations
 488 (Sec. 3.2). To isolate the role of $\langle \text{Ano} \rangle$ tokens, we
 489 generate $\langle \text{Diff} \rangle$ tokens using the unmodified, original
 490 visual features without salience adjustment.

491 As reported in Tab. 4 (*Effect of $\langle \text{Ano} \rangle$ Tokens*), removing
 492 $\langle \text{Ano} \rangle$ tokens consistently degrades performance across
 493 tasks. The drop reveals two important observations. First,
 494 anomaly-aware magnitude modulation enhances the
 495 expressiveness of the visual embeddings used for temporal
 496 differencing, allowing $\langle \text{Diff} \rangle$ tokens to capture clinically
 497 relevant cues rather than global appearance changes. Second,
 498 incorporating $\langle \text{Ano} \rangle$ tokens during inference adds
 499 complementary contextual information that stabilizes the
 500 interpretation of new findings. Together, these results show
 501 that $\langle \text{Ano} \rangle$ and $\langle \text{Diff} \rangle$ tokens form a mutually reinforcing
 502 pair: one grounds spatial anomaly cues, while the other
 503 captures temporal evolution, and both are required for
 504 accurate modeling of medical image progression.

Table 4. Effect of utilizing $\langle \text{Ano} \rangle$ tokens and comparison of visual feature extraction strategies.

	$\langle \text{Ano} \rangle$	Avg. F1 (\uparrow)	MMXU (\uparrow)
<i>Effect of $\langle \text{Ano} \rangle$ Tokens</i>			
$\langle \text{Diff} \rangle$ tokens only	×	–	0.635
MEDIC-AD	✓	–	0.655
<i>Feature Selection Strategy</i>			
Last layer	✓	90.9	0.619
Intermediate 4-layers	✓	91.6	0.635

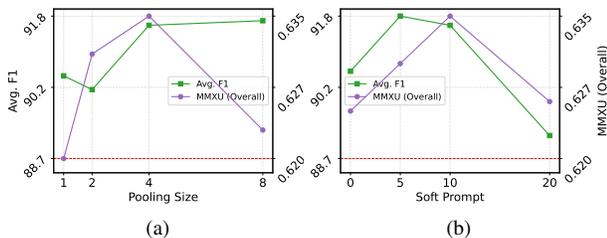


Figure 5. Hyperparameter sensitivity analysis on (a) query token pooling size and (b) visual soft prompt counts. The red line denotes the baseline performance of Lingshu [56].

5.2. Layer Selection in Visual Feature Extraction

In generating both $\langle \text{Ano} \rangle$ and $\langle \text{Diff} \rangle$ tokens, MEDIC-AD relies on visual features extracted from the vision encoder of the backbone VLM. To understand how visual representations influence anomaly and temporal token construction, we compare two configurations: using only the last-layer visual features, and aggregating intermediate-layer features together with the final representation.

As shown in Tab. 4 (*Feature Selection Strategy*), using intermediate-layer features consistently outperforms relying solely on the final hidden state. This advantage echoes prior findings [10, 23, 55, 60] that multi-level feature aggregation provides a broader range of semantic cues and improves downstream performance. In our setting, incorporating intermediate features not only enriches the anomaly representations but also produces feature embeddings that are more stable and informative for computing inter-image differences in temporal reasoning.

5.3. Impact of Hyperparameters

MEDIC-AD involves two key hyperparameters: (1) the number of generated $\langle \text{Ano} \rangle$ and $\langle \text{Diff} \rangle$ tokens, and (2) the number of soft prompts injected into the vision encoder. The number of $\langle \text{Ano} \rangle$ and $\langle \text{Diff} \rangle$ tokens is implicitly determined by the 2D pooling size applied to the magnitude-adjusted visual features used for token generation. Therefore, we investigate the effect of varying the pooling size. In parallel, the number of soft prompts influences how the extracted visual features are adapted to MEDIC-AD while also affecting the overall performance of

Table 5. Temporal image captioning performance on a real-world clinical dataset of 300 patients. Scores (0–1) indicate similarity to ground-truth descriptions, evaluated by GPT-4o [24].

Model	GPT-4o eval
Lingshu-7B	0.177
Medic-AD (Ours)	0.291

the backbone VLM, motivating a sensitivity analysis.

As illustrated in Fig. 5 (a), the model achieves consistently strong performance on both Anomaly Detection and MMXU benchmarks when the query-token pooling size is set to 4×4 , indicating that this granularity provides a favorable balance between spatial abstraction and anomaly localization. Similarly, the soft prompt sensitivity study in Fig. 5 (b) demonstrates that using 10 visual soft prompts yields the most stable and competitive results. We therefore adopt a pooling size of 4×4 and 10 soft prompts as the default configuration for MEDIC-AD.

5.4. Real-world Application

To validate the applicability of MEDIC-AD beyond public benchmarks, we conduct a real-world clinical study using chest X-ray pairs collected from 300 patients who visited a hospital for follow-up examinations. For each image pair, radiologists provide structured annotations describing the presence or absence of specific clinical findings, and the degree of change compared to the previous examination. Using this dataset, we formulate a temporal-difference captioning task in which the model must generate clinically consistent descriptions of symptom progression.

As shown in Tab. 5, MEDIC-AD outperforms Lingshu [56]—the strongest baseline in temporal reasoning (Tab. 2)—under GPT evaluation. These results indicate that MEDIC-AD remains effective not only in controlled benchmarks but also demonstrates robustness and reliability on real-world clinical data. Moreover, the model produces descriptions that align closely with expert assessments, highlighting its potential for integration into clinical workflows that demand interpretable and accurate temporal reasoning.

6. Conclusion

We introduced **MEDIC-AD**, a stage-wise medical VLM that strengthens clinical intelligence: lesion detection, temporal reasoning, and visual explainability through anomaly-aware and difference-token mechanisms. Our unified design enables the model to focus on abnormality cues, capture clinically meaningful changes between images, and provide accurate grounded visual evidence. Furthermore, evaluations on real-world hospital cases show robust alignment with expert assessments, indicating that MEDIC-AD offers a practical and reliable application for clinically usable medical VLMs.

577

References

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

- [1] Uwa O Aideyan, Kevin Berbaum, and Wilbur L Smith. Influence of prior radiologic information on the interpretation of radiographic examinations. *Academic Radiology*, 2(3): 205–208, 1995. 1
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *CVPR*, pages 2425–2433, 2015. 1
- [3] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 1, 2, 6
- [4] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022. 2, 6
- [5] Ujjwal Baid, Satyam Ghodasara, Suyash Mohan, Michel Bilello, Evan Calabrese, Errol Colak, Keyvan Farahani, Jayashree Kalpathy-Cramer, Felipe C Kitamura, Sarthak Pati, et al. The rsna-asnr-miccai brats 2021 benchmark on brain tumor segmentation and radiogenomic classification. *arXiv preprint arXiv:2107.02314*, 2021. 7
- [6] Shruthi Bannur, Stephanie Hyland, Qianchu Liu, Fernando Perez-Garcia, Maximilian Ilse, Daniel C Castro, Benedikt Boecking, Harshita Sharma, Kenza Bouzid, Anja Thieme, et al. Learning to exploit temporal structure for biomedical vision-language processing. In *CVPR*, pages 15016–15027, 2023. 2, 3
- [7] Jinan Bao, Hanshi Sun, Hanqiu Deng, Yinsheng He, Zhaoxiang Zhang, and Xingyu Li. Bmad: Benchmarks for medical anomaly detection. *arXiv preprint arXiv:2306.11876*, 2023. 3, 5, 7
- [8] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. The mvtec anomaly detection (mvtec ad) dataset: A comprehensive real-world dataset for unsupervised anomaly detection. Technical report, MVTEC Software GmbH, 2021. 3
- [9] Patrick Bilic, Patrick Christ, Hongwei Bran Li, Eugene Vorontsov, Avi Ben-Cohen, Georgios Kaissis, Adi Szeskin, Colin Jacobs, Gabriel Efrain Humpire Mamani, Gabriel Chartrand, et al. The liver tumor segmentation benchmark (lits). *MIA*, 84:102680, 2023. 7
- [10] Yunkang Cao, Jiangning Zhang, Luca Frittoli, Yuqi Cheng, Weiming Shen, and Giacomo Boracchi. Adaclip: Adapting clip with hybrid learnable prompts for zero-shot anomaly detection. *arXiv preprint arXiv:2407.15795*, 2024. 3, 8
- [11] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024. 1, 2, 6
- [12] Yeongjae Cho, Taehee Kim, Heejun Shin, Sungzoon Cho, and Dongmyung Shin. Pretraining vision-language model for difference visual question answering in longitudinal chest x-rays. *arXiv preprint arXiv:2402.08966*, 2024. 3
- [13] Muhammad E. H. Chowdhury, Tawsifur Rahman, Amith Khandakar, Rashid Mazhar, Muhammad Abdul Kadir, Zaid Bin Mahbub, Khandakar Reajul Islam, Muhammad Salman Khan, Atif Iqbal, Nasser Al Emadi, Mamun Bin Ibne Reaz, and Mohammad Tariqul Islam. Can ai help in screening viral and covid-19 pneumonia? *IEEE Access*, 8:132665–132676, 2020. 5, 6
- [14] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *NeurIPS*, 36: 49250–49267, 2023. 1, 2
- [15] Francesco Dalla Serra, Patrick Schrempf, Chaoyang Wang, Zaiqiao Meng, Fani Deligianni, and Alison Q. O’Neil. Grounding chest x-ray visual question answering with generated radiology reports. *arXiv preprint arXiv:2505.16624*, 2025. 3
- [16] Yossi Gandelsman, Alexei A. Efros, and Jacob Steinhardt. Interpreting clip’s image representation via text-based decomposition. *arXiv preprint arXiv:2310.05916*, 2023. 3
- [17] Zhaopeng Gu, Bingke Zhu, Guibo Zhu, Yingying Chen, Ming Tang, and Jinqiao Wang. Anomalygpt: Detecting industrial anomalies using large vision-language models. In *AAAI*, pages 1932–1940, 2024. 2, 3, 6
- [18] Ahmed Hamada. Br35h: Brain tumor detection 2020, 2020. 5, 6
- [19] Xuehai He, Yichen Zhang, Luntian Mou, Eric Xing, and Pengtao Xie. Pathvqa: 30000+ questions for medical visual question answering. *arXiv preprint arXiv:2003.10286*, 2020. 2, 5, 1
- [20] Junjie Hu, Yuanyuan Chen, and Zhang Yi. Automated segmentation of macular edema in oct using deep neural networks. *MIA*, 55:216–227, 2019. 7
- [21] Xinyue Hu, Lin Gu, Qiyuan An, Mengliang Zhang, Liangchen Liu, Kazuma Kobayashi, Tatsuya Harada, Ronald M. Summers, and Yingying Zhu. Expert knowledge-aware image difference graph representation learning for difference-aware medical visual question answering. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD ’23)*, pages 1–16. ACM, 2023. 3, 5
- [22] Yutao Hu, Tianbin Li, Quanfeng Lu, Wenqi Shao, Junjun He, Yu Qiao, and Ping Luo. Omnimedvqa: A new large-scale comprehensive evaluation benchmark for medical lvlm. In *CVPR*, pages 22170–22183, 2024. 2
- [23] Chaoqin Huang, Aofan Jiang, Jinghao Feng, Ya Zhang, Xinchao Wang, and Yanfeng Wang. Adapting visual-language models for generalizable anomaly detection in medical images. In *CVPR*, 2024. 3, 8
- [24] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 2, 3, 6, 8
- [25] Saahil Jain, Ashwin Agrawal, Adriel Saporta, Steven Q.H. Truong, Du Du Nguyen Duong, Tan Bui, Pierre Chambon, 633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689

- 690 Yuhao Zhang, Matthew P. Lungren, Andrew Y. Ng, Curtis P. Langlotz, and Pranav Rajpurkar. Radgraph: Extracting clinical entities and relations from radiology reports. *arXiv preprint arXiv:2106.14463*, 2021. 3
- 691
- 692
- 693
- 694 [26] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *ECCV*, pages 709–727. Springer, 2022. 4
- 695
- 696
- 697
- 698 [27] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. Pubmedqa: A dataset for biomedical research question answering. 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 2567–2577, 2019. 1
- 699
- 700
- 701
- 702
- 703
- 704
- 705 [28] Pranita Balaji Kanade and PP Gumaste. Brain tumor detection using mri images. *Brain*, 3(2):146–150, 2015. 5, 6
- 706
- 707 [29] Jaek Kim, Woojin Kim, Woohyeon Park, and Jaeyoung Do. Mmpb: It’s time for multi-modal personalization. *arXiv preprint arXiv:2509.22820*, 2025. 6
- 708
- 709
- 710 [30] Felipe Campos Kitamura. Head ct - hemorrhage, 2018. 5, 6
- 711 [31] Bennett Landman, Zhoubing Xu, J Igelsias, Martin Styner, T Langerak, and Arno Klein. Miccai multi-atlas labeling beyond the cranial vault—workshop and challenge. In *Proc. MICCAI Multi-Atlas Labeling Beyond Cranial Vault—Workshop Challenge*, page 12, 2015. 7
- 712
- 713
- 714
- 715 [32] 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. 2, 5, 1
- 716
- 717
- 718
- 719
- 720 [33] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. 1, 2
- 721
- 722
- 723
- 724 [34] 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. *NeurIPS*, 36:28541–28564, 2023. 1, 2, 6
- 725
- 726
- 727
- 728
- 729 [35] Junnan Li, Ramprasaath R. Selvaraju, Akhilesh D. Gotmare, Shafiq Joty, Caiming Xiong, and Steven C.H. Hoi. Align before fuse: Vision and language representation learning with momentum distillation. In *arXiv preprint arXiv:2107.07651*, 2021. 3
- 730
- 731
- 732
- 733
- 734 [36] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, pages 19730–19742. PMLR, 2023. 3
- 735
- 736
- 737
- 738 [37] Tang Li, Mengmeng Ma, and Xi Peng. Deal: Disentangle and localize concept-level explanations for vlms. *arXiv preprint arXiv:2407.14412*, 2024. 3
- 739
- 740
- 741 [38] Jie Lian, Jingyu Liu, Shu Zhang, Kai Gao, Xiaoqing Liu, Dingwen Zhang, and Yizhou Yu. A structure-aware relation network for thoracic diseases detection and segmentation. *IEEE Transactions on Medical Imaging*, 40(8):2042–2052, 2021. 3, 5, 7
- 742
- 743
- 744
- 745
- [39] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *NeurIPS*, 36:34892–34916, 2023. 1, 2, 3
- [40] Michael Moor, Qian Huang, Shirley Wu, Michihiro Yasunaga, Yash Dalmia, Jure Leskovec, Cyril Zarka, Eduardo Pontes Reis, and Pranav Rajpurkar. Med-flamingo: a multimodal medical few-shot learner. In *Machine Learning for Health (ML4H)*, pages 353–367. PMLR, 2023. 1
- [41] Linjie Mu, Zhongzhen Huang, Shengqian Qin, Yakun Zhu, Shaoting Zhang, and Xiaofan Zhang. Mmxu: A multi-modal and multi-x-ray understanding dataset for disease progression. *arXiv preprint arXiv:2502.11651*, 2025. 1, 2, 6
- [42] Xichen Pan, Li Dong, Shaohan Huang, Zhiliang Peng, Wenhu Chen, and Furu Wei. Kosmos-g: Generating images in context with multimodal large language models. *arXiv preprint arXiv:2310.02992*, 2023. 3
- [43] Yue Qiu, Shintaro Yamamoto, Kodai Nakashima, Ryota Suzuki, Kenji Iwata, Hirokatsu Kataoka, and Yutaka Satoh. Describing and localizing multiple changes with transformers. In *ICCV*, pages 1971–1980, 2021. 3
- [44] Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Eric Xing, Ming-Hsuan Yang, and Fahad S Khan. Glamm: Pixel grounding large multimodal model. In *CVPR*, pages 13009–13018, 2024. 1
- [45] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 5, 7
- [46] Andrew Sellergren, Sahar Kazemzadeh, Tiam Jaroensri, Atilla Kiraly, Madeleine Traverse, Timo Kohlberger, Shawn Xu, Fayaz Jamil, Cían Hughes, Charles Lau, et al. Medgemma technical report. *arXiv preprint arXiv:2507.05201*, 2025. 2
- [47] Omkar Thawkar, Abdelrahman Shaker, Sahal Shaji Mullappilly, Hisham Cholakkal, Rao Muhammad Anwer, Salman Khan, Jorma Laaksonen, and Fahad Shahbaz Khan. Xraygpt: Chest radiographs summarization using medical vision-language models. *arXiv preprint arXiv:2306.07971*, 2023. 1
- [48] Tao Tu, Shekoofeh Azizi, Danny Driess, Mike Schaeckermann, Mohamed Amin, Pi-Chuan Chang, Andrew Carroll, Charles Lau, Ryutaro Tanno, Ira Ktena, et al. Towards generalist biomedical ai. *Nejm Ai*, 1(3):AIoa2300138, 2024. 2, 5, 1
- [49] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In *CVPR*, pages 3156–3164, 2015. 1
- [50] Guoxin Wang, Jun Zhao, Xinyi Liu, Yanbo Liu, Xuyang Cao, Chao Li, Zhuoyun Liu, Qintian Sun, Fangru Zhou, Haoqiang Xing, et al. Citrus-v: Advancing medical foundation models with unified medical image grounding for clinical reasoning. *arXiv preprint arXiv:2509.19090*, 2025. 1, 2, 3, 5, 6, 7
- [51] Zifeng Wang, Zhenbang Wu, Dinesh Agarwal, and Jimeng Sun. Medclip: Contrastive learning from unpaired medical images and text. In *EMNLP*, page 3876, 2022. 2

- 804 [52] Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei
805 Chen, Zhuang Liu, In So Kweon, and Saining Xie. Con-
806 vnext v2: Co-designing and scaling convnets with masked
807 autoencoders. In *CVPR*, pages 16133–16142, 2023. 5
- 808 [53] Tinghao Xie, Xiangyu Qi, Yi Zeng, Yangsibo Huang,
809 Udari Madhushani Schwag, Kaixuan Huang, Luxi He, Boyi
810 Wei, Dacheng Li, Ying Sheng, et al. Sorry-bench: Systemat-
811 ically evaluating large language model safety refusal. *arXiv*
812 *preprint arXiv:2406.14598*, 2024. 6
- 813 [54] Yunfei Xie, Ce Zhou, Lang Gao, Juncheng Wu, Xianhang
814 Li, Hong-Yu Zhou, Sheng Liu, Lei Xing, James Zou, Cihang
815 Xie, et al. Medtrinity-25m: A large-scale multimodal dataset
816 with multigranular annotations for medicine. *arXiv preprint*
817 *arXiv:2408.02900*, 2024. 2
- 818 [55] Jiacong Xu, Shao-Yuan Lo, Bardia Safaei, Vishal M. Patel,
819 and Isht Dwivedi. Towards zero-shot anomaly detection and
820 reasoning with multimodal large language models. *arXiv*
821 *preprint arXiv:2502.07601*, 2025. 2, 3, 6, 8
- 822 [56] Weiwen Xu, Hou Pong Chan, Long Li, Mahani Aljunied,
823 Ruifeng Yuan, Jianyu Wang, Chenghao Xiao, Guizhen Chen,
824 Chaoqun Liu, Zhaodonghui Li, et al. Lingshu: A general-
825 ist foundation model for unified multimodal medical under-
826 standing and reasoning. *arXiv preprint arXiv:2506.07044*,
827 2025. 1, 2, 4, 6, 8
- 828 [57] Li Yang, Yan Xu, Chunfeng Yuan, Wei Liu, Bing Li, and
829 Weiming Hu. Improving visual grounding with visual-
830 linguistic verification and iterative reasoning. In *CVPR*,
831 pages 9499–9508, 2022. 1
- 832 [58] Linli Yao, Weiyang Wang, and Qin Jin. Image difference cap-
833 tioning with pre-training and contrastive learning. In *AAAI*,
834 pages 3108–3116, 2022. 3
- 835 [59] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi
836 Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming
837 Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline
838 multimodal understanding and reasoning benchmark for ex-
839 pert agi. In *CVPR*, pages 9556–9567, 2024. 5, 1
- 840 [60] Q Zhou and et al. Object-agnostic prompt learning for zero-
841 shot anomaly detection. *arXiv preprint arXiv:2310.18961*,
842 2023. 3, 8
- 843 [61] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mo-
844 hamed Elhoseiny. Minigt-4: Enhancing vision-language
845 understanding with advanced large language models. *arXiv*
846 *preprint arXiv:2304.10592*, 2023. 1, 2
- 847 [62] Yang Zou, Jongheon Jeong, Latha Pemula, Dongqing Zhang,
848 and Onkar Dabeer. Spot-the-difference self-supervised pre-
849 training for anomaly detection and segmentation. In *ECCV*,
850 pages 392–408, 2022. 3