

A MEDICAL MULTIMODAL DIAGNOSTIC FRAMEWORK INTEGRATING VISION-LANGUAGE MODELS AND LOGIC TREE REASONING

Zelin Zang^{1,2}, Wenyi Gu¹, Siqi Ma², Dan Yang³, Yue Shen³, Zhu Zhang⁴, Guohui Fan⁴, Wing-Kuen Ling¹, Fuji Yang¹

¹ Tsientang Institute of Advanced Study (TIAS), Hangzhou, China

² Westlake University, Hangzhou, China ³ Ant Group, Hangzhou, China

⁴ China-Japan Friendship Hospital, Beijing, China zangzelin@westlake.edu.cn,

ABSTRACT

With the rapid growth of large language models (LLMs) and vision-language models (VLMs) in medicine, simply integrating clinical text and medical imaging does not guarantee reliable reasoning. Existing multimodal models often produce hallucinations or inconsistent chains of thought, limiting clinical trust. We propose a diagnostic framework built upon LLaVA that combines vision-language alignment with logic-regularized reasoning. The system includes an input encoder for text and images, a projection module for cross-modal alignment, a reasoning controller that decomposes diagnostic tasks into steps, and a logic tree generator that assembles stepwise premises into verifiable conclusions. Evaluations on MedXpertQA and other benchmarks show that our method improves diagnostic accuracy and yields more interpretable reasoning traces on multimodal tasks, while remaining competitive on text-only settings. These results suggest a promising step toward trustworthy multimodal medical AI.

Index Terms— Medical Multimodal Diagnosis; Vision-Language Model; Logic Tree Reasoning; Explainable Artificial Intelligence;

1. INTRODUCTION

Deep learning has greatly advanced medical AI. Large language models (LLMs) [1] and vision-language models (VLMs) [2] can now jointly process clinical text, history, and images, achieving promising results across various medical domains including disease classification [3] and clinical prediction [4]. This often improves predictions — for example, in our tests, models combining CT findings with symptom descriptions narrowed differential diagnoses better than text-only systems. Yet simply adding modalities does not guarantee sound reasoning. Benchmarks such as MedXpertQA [5], VQA-RAD [6], PathVQA [7], and PubMedQA [8] show that multimodal models may contradict themselves or ignore key evidence. Such inconsistencies limit trust and slow adoption, since physicians must verify how conclusions are reached [9].

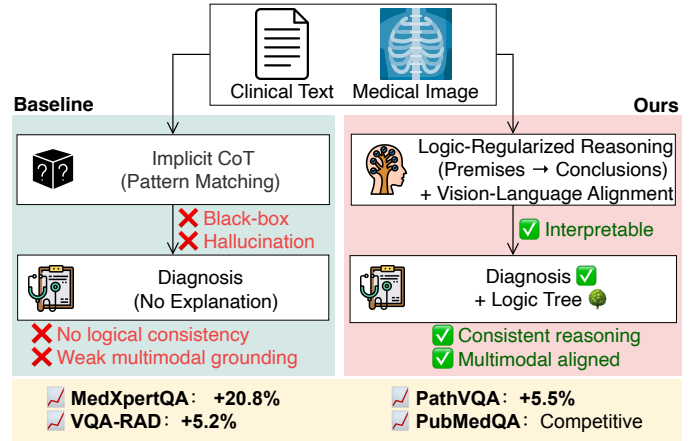


Fig. 1. Overview. Left: baseline VLMs rely on implicit CoT, causing hallucinations, inconsistency, and weak grounding. Right: our framework integrates vision-language alignment and logic-regularized reasoning, yielding traceable logic trees and consistent diagnoses.

Despite progress, strong VLMs still give confident but wrong answers, or hallucinations [10], especially when text and image evidence conflict. In pilot runs, some even contradicted radiology reports by over-relying on text — dangerous in clinical practice, where unsupported predictions can mislead physicians.

On closer inspection, these errors seemed to stem from implicit pattern matching rather than structured reasoning, which made their outputs hard to audit. To address this, we tried making the reasoning process explicit instead of directly predicting a label. Early attempts were unstable, but adding formal logic constraints stabilized training and surprisingly improved both accuracy and interpretability. Our current model now learns to arrange premises step by step before reaching a conclusion — similar to how clinicians justify decisions in tumor boards — producing reasoning chains that are easier for physicians to inspect and question.

Building on LLaVA [11], which has demonstrated strong vision-language alignment capabilities and efficient multi-

modal fusion, we designed a system with four cooperating parts: an input encoder for clinical text and CT/MRI images, a vision-language alignment module that maps image features into the language space, a reasoning controller that breaks down diagnostic tasks into intermediate steps, and a logic tree generator that assembles those steps into a verifiable premise-conclusion chain. In early trials, we noticed that logic rules alone had limited impact — the model still ignored subtle visual cues. Multimodal alignment turned out to be essential: once we projected visual features into the same space as text, the reasoning controller produced much more consistent outputs. Our final design combines formal logic constraints with vision-language alignment, yielding reasoning trees that physicians can check. Experiments on MedXpertQA [5] and other benchmarks confirmed the benefit: accuracy improved, and the reasoning traces were easier to verify, which we see as a step toward more trustworthy medical AI.

2. RELATED WORK

Vision–Language Models in Healthcare. Recent multimodal studies have produced several specialized VLMs for medicine. CheXzero showed that paired image-text pretraining could reach radiologist-level performance on chest X-rays without manual labels, hinting that large unannotated corpora might replace expensive annotation. Med-Flamingo [12] adapted OpenFlamingo to the clinical setting and reported a 20 % improvement in blinded physician ratings when rationales were provided — a result that convinced us that explanation quality really matters in practice. LLaVA-Med [13] further demonstrated that instruction-tuning on PubMed images with GPT-4 captions yields strong biomedical VQA performance. More recent systems such as UMIT [14] and HealthGPT [15] attempt to unify multiple imaging tasks and support both comprehension and generation. Beyond diagnostic imaging, recent work has explored multimodal integration in other biomedical contexts such as spatial transcriptomics [16] and single-cell analysis [17], demonstrating the broader applicability of vision-language methods. These efforts show impressive gains in recognition, but we found that most still lack structured reasoning that clinicians can audit. **Logic-Based Reasoning in Medical LLMs.** Explicit reasoning has been explored to increase trust in model outputs. Med-PaLM [1] and Med-PaLM 2 [18] fine-tuned general LLMs to produce step-by-step chains of thought, though these often remain opaque and difficult to verify. MedLA [19] goes further by using multi-agent dialogue to refine logic trees, while MDAgents [20] coordinates multiple expert agents under a moderator. These approaches improved clinical rigor but came with heavy computational overhead. We chose a simpler path: embedding logic constraints directly into a single-model chain-of-thought, aiming to keep the reasoning consistent and interpretable without requiring

multi-agent orchestration.

3. METHOD

Overall Framework. We extend LLaVA to the clinical setting by adding explicit logic regularization. Our system proceeds in three stages, illustrated in Figure 2. First, we embed medical images and clinical narratives into a shared representation space so that visual and textual cues can interact early. Next, the model generates several candidate chains-of-thought, which we refine with a dynamic optimization scheme (DAPO) [21] that balances three signals: diagnostic accuracy, logical consistency, and image-text grounding. Finally, we parse the best reasoning path into syllogistic triads and assemble them into a logic tree that clinicians can inspect step by step. During development, we tried freezing the vision encoder to save compute, but convergence became unstable and the model ignored subtle imaging findings. Jointly training all components solved this problem, so we adopted a fully end-to-end optimization strategy.

Visual & Text Encoder. We adopt a ViT [22] backbone pretrained on vision–language tasks as the image encoder. Each 2D medical image is split into patches and converted into visual tokens v_1, \dots, v_M . For volumetric scans (CT or MRI), we process each slice with the same encoder and then use a slice-fusion transformer with multi-head attention to combine information across slices. In early tests, we tried simply averaging slice features, but this blurred subtle lesions and hurt performance on multi-slice CT cases, so we switched to the attention-based fusion module. Clinical text T —including patient history, imaging findings, and diagnostic questions—is tokenized and embedded using a pretrained LLM backbone (LLaMA [23] or Vicuna [24]). We project the resulting textual tokens t_1, \dots, t_N into the same hidden space as the visual tokens. This early projection allows the subsequent attention layers to reason jointly over text and image features rather than treating them as separate streams.

Vision–Language Alignment. To enable joint reasoning, we project visual tokens into the LLM hidden space through a learned projection matrix $W_{\text{proj}} \in \mathbb{R}^{d_h \times d_v}$:

$$h_i^{\text{vis}} = W_{\text{proj}} v_i. \quad (1)$$

The projected features are interleaved with textual embeddings h_j^{txt} before entering the LLaVA backbone, allowing subsequent self-attention layers to fuse cross-modal context. During early trials, we observed that without explicit alignment, the model often ignored subtle imaging cues. To address this, we compute global embeddings z_v and z_t (mean-pooled over tokens) and apply a CLIP-style [25] InfoNCE loss, which has been shown effective for aligning noisy multimodal medical data [26]:

$$\mathcal{L}_{\text{align}} = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(S(z_{v_i}, z_{t_i})/\tau)}{\sum_{j=1}^B \exp(S(z_{v_i}, z_{t_j})/\tau)}, \quad (2)$$

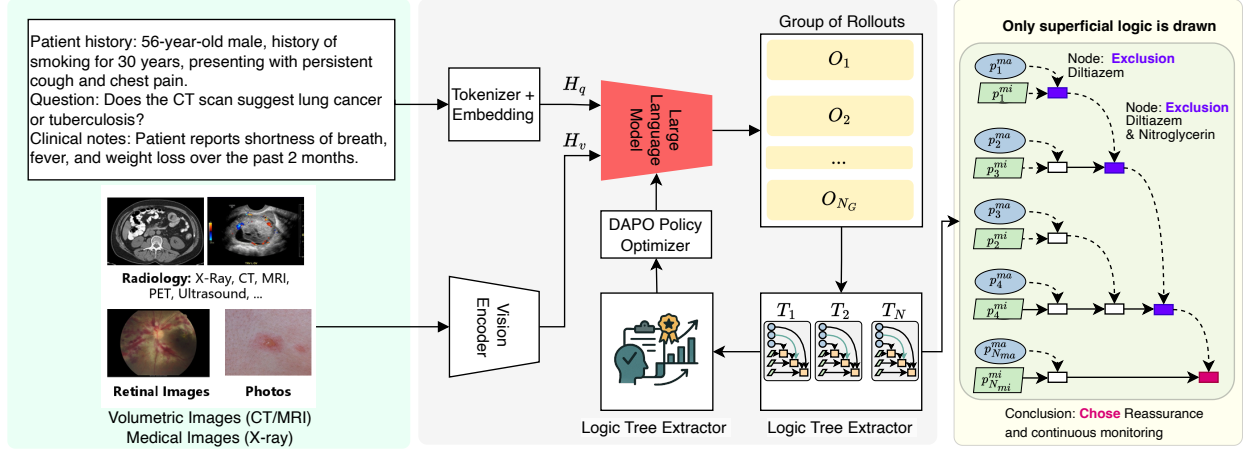


Fig. 2. Multimodal logic-regularized reasoning framework. Patient history, clinical notes, and medical images (e.g., CT, MRI, X-ray) are encoded by text and vision encoders. Their embeddings are fed into a large language model with a DAPO policy optimizer, producing multiple candidate reasoning rollouts. Each rollout is parsed by the logic tree extractor into syllogistic premises and conclusions, forming verifiable logic trees. The final diagnosis is obtained together with a traceable reasoning chain, improving both accuracy and interpretability.

where $S(z_v, z_t) = \frac{z_v^\top z_t}{\|z_v\| \|z_t\|}$ is cosine similarity, τ is a temperature, and B is batch size. Recent work has shown that soft contrastive objectives [27] and diffusion-based augmentation [28] can further enhance such alignment, though we leave these extensions for future work. This additional loss encouraged better grounding of visual tokens and reduced hallucination in ablation studies.

Prompt-Based Reasoning with Rollouts. We elicit explicit chains-of-thought (CoTs) [29] by combining instructional prompts (e.g., “Let’s reason step by step”) with a small set of premise–conclusion exemplars. The model produces multiple rollouts O_k , each a trajectory of K reasoning steps. Rather than relying solely on likelihood maximization, we introduce a logic-based regularizer:

$$\mathcal{L}_{\text{logic}} = \frac{1}{K} \sum_{k=1}^K \left(1 - f_{\text{logic}}(p_{\text{maj}}^{(k)}, p_{\text{min}}^{(k)}, c^{(k)}) \right), \quad (3)$$

where $(p_{\text{maj}}^{(k)}, p_{\text{min}}^{(k)}) \rightarrow c^{(k)}$ is a syllogistic triad and $f_{\text{logic}} \in [0, 1]$ is a rule-based verifier that checks entailment, flags contradictions, and penalizes unsupported conclusions. Specifically, f_{logic} assigns 1.0 when the conclusion validly follows from the premises via modus ponens or modus tollens, 0.5 for weak inferences, and 0.0 for contradictions or non-sequiturs. We use Dynamic Advantage Policy Optimization (DAPO) [21], a variant of Proximal Policy Optimization (PPO), which reweights advantages by combining three signals: diagnostic correctness, logic consistency, and vision–language grounding. In practice, we found that this multi-objective training stabilized learning and improved both accuracy and explanation quality.

Logic Tree Generator. The final stage parses each reasoning trajectory into a logic tree [19], where every edge rep-

resents a syllogistic triad. This structured representation allows us to inspect intermediate conclusions and trace back incorrect predictions during error analysis. Given a multimodal input $x = (I, T)$, the tree outputs a final diagnosis \hat{y} , and the model is trained with a cross-entropy loss:

$$\mathcal{L}_{\text{diag}} = - \sum_{c=1}^C \mathbf{1}[y = c] \log p_{\theta}(c | x), \quad (4)$$

where $y \in \{1, \dots, C\}$ is the ground-truth label and $p_{\theta}(c | x)$ the predicted class probability. During development, we experimented with margin-based objectives but found that standard cross-entropy yielded more stable convergence when combined with logic regularization.

Training Objective. Our final objective integrates three components: diagnostic accuracy, logical consistency, and multimodal alignment:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{diag}} + \lambda_{\text{logic}} \mathcal{L}_{\text{logic}} + \lambda_{\text{align}} \mathcal{L}_{\text{align}}. \quad (5)$$

The weights λ_{logic} and λ_{align} are tuned on the validation set to balance prediction performance and reasoning quality. In practice, we found that setting λ_{logic} too high caused the model to overfit on rule satisfaction at the expense of accuracy, so we adopt a moderate weighting that preserves both interpretability and clinical relevance.

4. EXPERIMENTS

Datasets. We evaluate on four complementary QA/VQA benchmarks spanning expert-level reasoning, radiology, pathology, and biomedical text inference. Our *primary benchmark* is *MedXpertQA* [5], with 4,460 expert-level questions across 17 specialties and 11 organ systems, including

Table 1. Results on VQA-RAD, PathVQA, and PubMedQA. Accuracy (%) and explanation quality (ROUGE-L, higher is better). Our model outperforms baselines on multimodal tasks and remains competitive on PubMedQA.

| Model | VQA-RAD | | PathVQA | | PubMedQA | |
|-------------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | Acc | R-L | Acc | R-L | Acc | R-L |
| Med-PaLM 2 [18] | 63.4 | 35.1 | <u>60.2</u> | <u>33.8</u> | 79.6 | 44.2 |
| GPT-4V | 65.1 | <u>34.7</u> | 58.6 | 32.1 | 76.9 | 42.8 |
| BioMedCLIP [30] | 61.2 | 31.6 | 56.1 | 30.2 | 70.5 | 39.7 |
| Med-Flamingo [12] | 66.8 | 33.9 | 59.3 | 31.0 | 72.4 | 40.5 |
| LLaVA-Med [13] | 62.5 | 31.1 | 55.4 | 29.7 | 69.1 | 38.9 |
| MedRAG [31] | <u>67.2</u> | 34.1 | 60.0 | 31.4 | 77.3 | 43.1 |
| Ours | 72.4 | 38.5 | 65.7 | 36.2 | <u>78.8</u> | <u>43.9</u> |

2,005 multimodal cases paired with 2,839 images (radiology, pathology, clinical photos, charts). Each case requires integrating text and images for diagnostic or treatment reasoning. We follow the official train/validation/test split to avoid patient overlap. Since training only on MedXpertQA risked overfitting to question style, we added three datasets for generalization: (1) *VQA-RAD* [6], 315 radiology images and 3,515 Q–A pairs (CT, MRI, X-ray); (2) *PathVQA* [7], 4,998 pathology images and 32,799 Q–A pairs; (3) *PubMedQA* [8], 211k text-only pairs from PubMed abstracts.

Model Variants and Evaluation Metrics. We compare our full *Logic-Regularized VLA*, combining vision–language alignment with logic-regularized reasoning, against a *No-Logic* ablation where logical constraints are removed. This isolates the effect of logic regularization on accuracy and explanation quality. Pilot tests also showed removing vision inputs greatly reduced accuracy, confirming image grounding is essential. We evaluate with *diagnostic accuracy* and *ROUGE-L*. Accuracy is computed after VQA normalization (lowercasing, punctuation removal, synonym mapping), while ROUGE-L measures reasoning coverage and coherence by comparing generated chains-of-thought against reference explanations. Although ROUGE-L primarily captures lexical overlap, prior work [29] has shown it correlates well with human judgments of reasoning quality in medical contexts where step-by-step explanations follow similar logical structures. We report mean \pm std over three seeds, using McNemar’s test for accuracy and paired bootstrap for ROUGE-L. Models with similar accuracy often differ in ROUGE-L, highlighting the need for both metrics.

Main Results. Tables 1 and 2 summarize results across four benchmarks. Our model achieves the best overall performance on multimodal tasks, with particularly gains on tree benchmarks. On the text-only PubMedQA, performance is on par with specialized text-only systems, suggesting that the added visual components do not harm purely textual reasoning. Interestingly, we found that gains were largest on cases requiring integration of subtle imaging findings with clinical

Table 2. Results on MedXpertQA. We report *Diagnostic Accuracy* (%) and *Explanation Quality* (ROUGE-L, higher is better). Our logic-regularized vision–language model substantially outperforms all baselines.

| Model | Accuracy (%) | ROUGE-L |
|------------------------|--------------|-------------|
| o1 | 56.3 | 41.2 |
| GPT-4V | 42.8 | 34.5 |
| Claude-3.5-Sonnet [32] | 33.2 | 32.0 |
| Gemini-1.5-Pro [33] | 34.1 | 32.5 |
| QVQ-72B-Preview [34] | 33.6 | 32.8 |
| Qwen2.5-VL-72B [35] | 30.0 | 31.0 |
| Ours | 77.1 | 41.6 |

history, which aligns with our design goal of improving multimodal consistency. The ablation study (Table 3) reveals that DAPO contributes 3.9% accuracy gain over standard PPO, likely because its dynamic advantage reweighting balances hard multimodal cases with straightforward text-only examples, preventing the model from collapsing to easier sub-tasks during optimization. **Analysis.** Qualitative review by two domain experts indicates that our model produces more coherent, stepwise explanations compared to baselines. Vision–language baselines often identify the correct image region but fail to link it explicitly to the diagnosis, resulting in “black-box” predictions. In contrast, logic regularization encourages the model to generate premise–conclusion chains, which not only improve interpretability but also helped reviewers identify when the model made an incorrect inference. We also observed that some failure cases stem from missing clinical context rather than model errors, highlighting opportunities for future integration with retrieval-based methods.

Ablation Study. We ran a series of ablations by selectively removing vision input, logic regularization, alignment loss, and DAPO optimization. As shown in Table 3, vision input had the biggest impact, confirming that image grounding is crucial for clinical reasoning. Logic regularization mainly improved the coherence of reasoning chains, while alignment loss helped reduce hallucinations. Removing DAPO made training noticeably less stable across runs. Together, these results indicate that all components contribute to final performance, with logic and alignment playing complementary roles in reliability.

5. CONCLUSION

In this work, we proposed a logic-regularized multimodal diagnostic framework that integrates visual–language alignment with structured reasoning. Experiments on four medical QA/VQA benchmarks show consistent improvements in both accuracy and explanation quality, with logical trees offering transparent reasoning paths for clinical review. While these results on benchmark datasets are promising, real-world clin-

Table 3. Ablation study on MedXpertQA. V: Vision, L: Logic loss, A: Alignment loss, D: DAPO. Each component contributes to performance: V improves accuracy, L enhances reasoning quality, A reduces hallucination, and D stabilizes optimization.

| Model | V | L | A | D | Acc. (%) | R-L |
|------------------|---|---|---|---|-------------|-------------|
| Full Model | ✓ | ✓ | ✓ | ✓ | 77.1 | 41.6 |
| - Logic Loss | ✓ | × | ✓ | ✓ | 72.3 | 35.7 |
| - Alignment Loss | ✓ | ✓ | × | ✓ | 70.7 | 39.1 |
| - DAPO | ✓ | ✓ | ✓ | × | 73.2 | 37.6 |
| - Vision | × | ✓ | ✓ | ✓ | 52.0 | 33.9 |

ical deployment would require further validation in actual practice settings. Nonetheless, our framework demonstrates a meaningful step toward more interpretable and reliable multimodal medical reasoning, offering both improved diagnostic accuracy and clearer reasoning chains that can facilitate physician verification. Future work could explore incorporating interpretable dimensionality reduction techniques [36] to better visualize the learned reasoning structures and enhance clinical understanding. **Compliance with Ethical Standards.** The authors used LLM to assist with language editing and polishing. All technical content, results, and conclusions are solely authored and verified by the authors.

6. REFERENCES

- [1] Karan Singhal, Shekoofeh Azizi, Towaki Tu, and et al., “Large language models encode clinical knowledge,” *Nature*, vol. 620, no. 7972, pp. 172–180, 2023.
- [2] Ekin Tiu, Ellie Talus, , and Pranav Rajpurkar, “Expert-level detection of pathologies from unannotated chest x-ray images via self-supervised learning,” *Nature Biomedical Engineering*, 2022.
- [3] Yaoting Sun, Sathiyamoorthy Selvarajan, Zelin Zang, Wei Liu, Yi Zhu, Hao Zhang, Wanyuan Chen, Hao Chen, Lu Li, Xue Cai, et al., “Artificial intelligence defines protein-based classification of thyroid nodules,” *Cell Discovery*, vol. 8, no. 1, pp. 85, 2022.
- [4] Kai Zhou, Yaoting Sun, Lu Li, Zelin Zang, Jing Wang, Jun Li, Junbo Liang, Fangfei Zhang, Qiushi Zhang, Weigang Ge, et al., “Eleven routine clinical features predict covid-19 severity uncovered by machine learning of longitudinal measurements,” *Computational and Structural Biotechnology Journal*, vol. 19, pp. 3640–3649, 2021.
- [5] Yuxin Zuo, Shang Qu, Yifei Li, Zhangren Chen, Xuekai Zhu, Ermo Hua, Kaiyan Zhang, Ning Ding, and Bowen Zhou, “Medxpertqa: Benchmarking expert-level medical reasoning and understanding,” *arXiv*, 2025.
- [6] 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*, vol. 5, no. 1, pp. 1–10, 2018.
- [7] Xuehai He, Yichen Zhang, Luntian Mou, Eric Xing, and Pengtao Xie, “Pathvqa: 30000+ questions for medical visual question answering,” *arXiv*, 2020.
- [8] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu, “Pubmedqa: A dataset for biomedical research question answering,” in *EMNLP-IJCNLP*, 2019, pp. 2567–2577.
- [9] Jean-Christophe Bélisle-Pipon, “Why we need to be careful with llms in medicine,” *Frontiers in Medicine*, vol. 11, pp. 1495582, 2024.
- [10] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung, “Survey of hallucination in natural language generation,” *ACM Computing Surveys*, vol. 55, no. 12, pp. 248:1–248:38, 2023.
- [11] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee, “Visual instruction tuning,” *NeurIPS*, 2023.

- [12] Michael Moor, Qian Huang, Shirley Wu, Michihiro Yasunaga, Cyril Zakka, Yash Dalmia, Eduardo Pontes Reis, Pranav Rajpurkar, and Jure Leskovec, “Med-flamingo: A multimodal medical few-shot learner,” *arXiv*, 2023.
- [13] Chunyuan Li and Jianfeng Gao, “Llava-med: Training a large language-and-vision assistant for biomedicine in one day,” *arXiv preprint arXiv:2306.00890*, 2023.
- [14] H. Yu, W. Zhang, J. Zhao, et al., “Umit: Unifying medical imaging tasks via vision-language models,” *arXiv*, 2025.
- [15] Tianwei Lin and Beng Chin Ooi, “Healthgpt: A medical large vision-language model,” 2025.
- [16] Zelin Zang, Liangyu Li, Yongjie Xu, Chenrui Duan, Yue Shen, Yi Sun, Zhen Lei, and Stan Z Li, “Must: multiple-modality structure transformation for single-cell spatial transcriptomics,” *Briefings in Bioinformatics*, vol. 26, no. 4, pp. bbaf405, 2025.
- [17] Yongjie Xu, Zelin Zang, Bozhen Hu, Yue Yuan, Cheng Tan, Jun Xia, and Stan Z Li, “Complex hierarchical structures analysis in single-cell data with poincaré deep manifold transformation,” *Briefings in Bioinformatics*, vol. 26, no. 1, pp. bbae687, 2025.
- [18] Karan Singhal and Others, “Toward expert-level medical question answering with large language models,” *Nature Medicine*, 2025.
- [19] Siqi Ma and Zelin Zang, “MedLA: A Logic-Driven Multi-Agent Framework for Complex Medical Reasoning with Large Language Models,” *Under review*, 2024.
- [20] Doeun Kim, Xiaoxiao Wang, et al., “Mdagents: An adaptive collaboration of llms for medical decision-making,” *arXiv preprint arXiv:2406.06782*, 2024.
- [21] Qiyang Yu and Mingxuan Wang, “Dapo: An open-source llm reinforcement learning system at scale,” 2025.
- [22] Alexey Dosovitskiy and Neil Houlsby, “An image is worth 16x16 words: Transformers for image recognition at scale,” *ICLR*, 2021.
- [23] Hugo Touvron and Guillaume Lample, “Llama: Open and efficient foundation language models,” *arXiv*, 2023.
- [24] Wei-Lin Chiang and Eric P. Xing, “Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality,” March 2023.
- [25] Alec Radford and Ilya Sutskever, “Learning transferable visual models from natural language supervision,” in *ICML*. 2021, pp. 8748–8763, PMLR.
- [26] Zijia Song, Zelin Zang, Yelin Wang, Guozheng Yang, Jiangbin Zheng, Wanyu Chen, and Stan Z Li, “Gentle-clip: Exploring aligned semantic in low-quality multimodal data with soft alignment,” *arXiv preprint arXiv:2406.05766*, 2024.
- [27] Zelin Zang, Lei Shang, Senqiao Yang, Fei Wang, Baigui Sun, Xuansong Xie, and Stan Z Li, “Boosting novel category discovery over domains with soft contrastive learning and all in one classifier,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023, pp. 11858–11867.
- [28] Zelin Zang, Hao Luo, Kai Wang, Panpan Zhang, Fan Wang, Stan Li, and Yang You, “Diffaug: Enhance unsupervised contrastive learning with domain-knowledge-free diffusion-based data augmentation,” in *International Conference on Machine Learning (ICML)*, 2024.
- [29] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc V. Le, and Denny Zhou, “Chain-of-thought prompting elicits reasoning in large language models,” in *NeurIPS*, 2022.
- [30] Sheng Zhang and Hoifung Poon, “Biomedclip: A multimodal biomedical foundation model pretrained from fifteen million scientific image-text pairs,” *arXiv*, 2023.
- [31] Guangzhi Xiong, Qiao Jin, Xiao Wang, Minjia Zhang, Zhiyong Lu, and Aidong Zhang, “Improving retrieval-augmented generation in medicine with iterative follow-up questions,” in *Biocomputing*, 2024.
- [32] Anthropic, “Claude-3.5 sonnet model card addendum,” Model Card / Addendum, 2024.
- [33] Google DeepMind and Google AI, “Gemini 1.5: Unlocking multimodal understanding across models including gemini-1.5 pro,” *Google Technical Report*, 2024.
- [34] Qwen Team, “Qvq-72b preview: A visual reasoning model by qwen,” Blog / Model Preview, 2024.
- [35] S. Bai et al., “Qwen2.5-vl-72b technical report,” *arXiv*, 2025.
- [36] Zelin Zang, Shenghui Cheng, Hanchen Xia, Liangyu Li, Yaoting Sun, Yongjie Xu, Lei Shang, Baigui Sun, and Stan Z Li, “Dmt-ev: An explainable deep network for dimension reduction,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 30, no. 3, pp. 1710–1727, 2024.