

Has Multimodal Learning Delivered Universal Intelligence in Healthcare? A Comprehensive Survey

Qika Lin, *Member, IEEE*, Yifan Zhu, *Member, IEEE*, Xin Mei, Ling Huang, Jingying Ma, Kai He, *Member, IEEE*, Zhen Peng, Erik Cambria, *Fellow, IEEE*, Mengling Feng, *Senior Member, IEEE*

Abstract—The rapid development of artificial intelligence has constantly reshaped the field of intelligent healthcare and medicine. As a vital technology, multimodal learning has increasingly garnered interest due to data complementarity, comprehensive modeling form, and great application potential. Currently, numerous researchers are dedicating their attention to this field, conducting extensive studies and constructing abundant intelligent systems. Naturally, an open question arises that *has multimodal learning delivered universal intelligence in healthcare?* To answer the question, we adopt three unique viewpoints for a holistic analysis. Firstly, we conduct a comprehensive survey of the current progress of medical multimodal learning from the perspectives of datasets, task-oriented methods, and universal foundation models. Based on them, we further discuss the proposed question from five issues to explore the real impacts of advanced techniques in healthcare, from data and technologies to performance and ethics. The answer is that current technologies have **NOT** achieved universal intelligence and there remains a significant journey to undertake. Finally, in light of the above reviews and discussions, we point out ten potential directions for exploration towards the goal of universal intelligence in healthcare.

Index Terms—Intelligent healthcare, medical intelligence, multimodal learning, foundation model, medical vision-language.

1 INTRODUCTION

RECENT years have seen the remarkable progress of artificial intelligence (AI) across the healthcare and medicine domain [1]. AI techniques have demonstrated substantial potential in various medical scenarios, including medical imaging analysis [2], disease diagnosis [3], drug discovery [4], personalized treatment [5], and medical QA (question-answering) [6], aiming to provide automated and customized expert-level advice or recommendations to alleviate the burden on both patients and physicians. Nonetheless, these studies or applications typically consider only single-modality data, *e.g.*, medical image or text, which could result in diminished performance and may not accurately represent authentic application scenarios [7].

As the healthcare domain ceaselessly produces an increasing volume and variety of data, ranging from medical images and clinical notes to genomic profiles and biosensor readings, the need for effective multimodal learning approaches becomes paramount [7], [8], [9], [10]. On the one hand, multimodal AI models, capable of integrating and learning from these heterogeneous data streams, hold the promise of unlocking a comprehensive and nuanced understanding of complex medical phenomena. By capturing complementary semantic information [11] (as shown in

Figure 1) and intricate relationships across different modalities [12], these models provide clinicians with a holistic view of patients' conditions, enabling more proactive monitoring, accurate diagnoses, and personalized treatment plans. On the other hand, multimodal learning further broadens the application prospects of intelligent models in the healthcare field. For instance, if a patient needs to ask about their skin condition, articulating it verbally (*e.g.*, using conventional language QA systems) can be challenging. A visual question-answering (VQA) system becomes incredibly useful, as it can combine intuitive images uploaded by patients to make more accurate and comprehensive diagnoses. Given the significant research importance and application value of multimodal healthcare, recent years have witnessed an extensive amount of research dedicated to this particular subject, with a clear rising trend. The advancement in technologies has progressed from utilizing specific models, such as convolutional neural network (CNN) [13], recurrent neural network (RNN) [14], graph neural network (GNN) [15], and Transformer [16], to the adoption of a strategy involving pre-training and fine-tuning. The latter has emerged as the prevailing focus and trending topic, which is inspired by the powerful foundational models (FMs) in the general domain, like CLIP [17], ChatGPT¹, GPT-4², and multimodal large language model (MLLM) [18]. These studies have made significant advancements in numerous tasks of multimodal healthcare, *e.g.*, multi-modality image fusion, report generation (RG), VQA, cross-modal retrieval, text-augmented image processing, and cross-modal image generation. The evolution has ultimately led to the development of FMs

- It is updating online at: <https://github.com/DeepReasoning/aihealth>.
- Qika Lin, Ling Huang, Jingying Ma, Kai He, and Mengling Feng are with the Saw Swee Hock School of Public Health, National University of Singapore, 117549, Singapore.
- Yifan Zhu is with the School of Computer Science, Beijing University of Posts and Telecommunications, Beijing, 100876, China.
- Xin Mei is with the School of Automation, Northwestern Polytechnical University, Xi'an, China.
- Zhen Peng is with the School of Computer Science and Technology, Xi'an Jiaotong University, Xi'an, Shaanxi 710049, China.

1. <https://openai.com/chatgpt>
2. <https://openai.com/gpt-4>

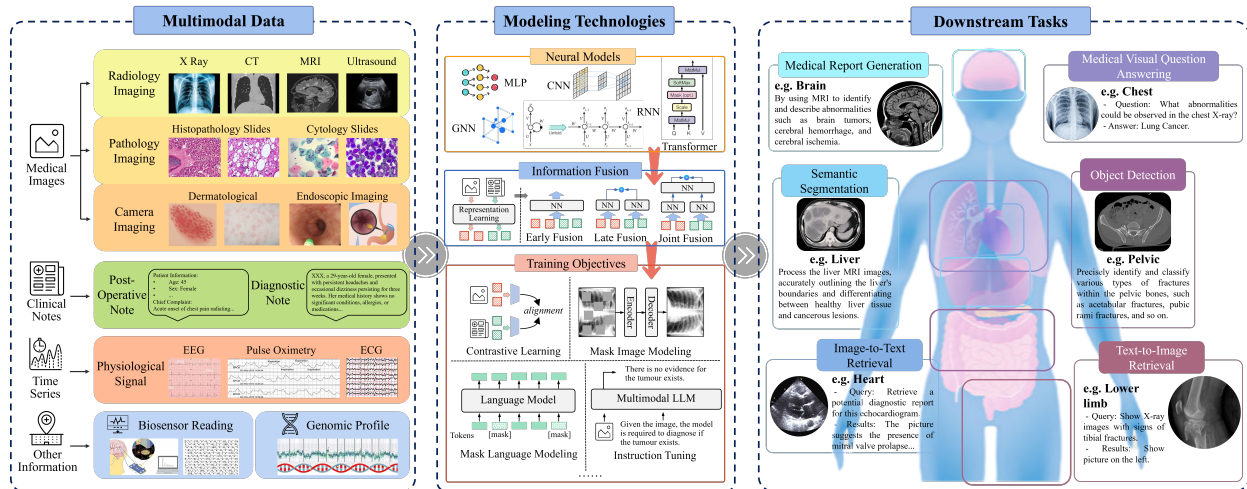


Fig. 1: Illustration of multimodal learning in healthcare, from perspectives of multimodal data and modeling technologies to downstream tasks.

capable of handling various medical tasks. We summarize the overall architecture in Figure 1.

Despite these seemingly enormous achievements, it is unclear how far existing research has progressed. More importantly, the trust among doctors and patients in applying existing methods to real-world scenarios is a significant question [7], [19]. To this end, we carry out this survey to answer the following open research question: *has multimodal learning delivered universal intelligence in healthcare?*, i.e., has multimodal learning delivered an advanced AI healthcare form with a broad range of cognitive abilities, comprehension of various situations, and practical applications? By answering this question, we aim to provide researchers with a comprehensive global overview of the advancements made towards achieving this objective, the remaining challenges to be addressed, and the necessary steps to be taken. To achieve this, we conduct our analysis from the following three dimensions: 1) We first comprehensively review the current progress of medical multimodal learning from the perspectives of datasets, task-oriented techniques, and universal FMs. 2) We discuss the open question from five issues to explore the real impacts of current advanced techniques in healthcare applications, from data and technologies to performance and ethics. We find that current technologies have **NOT** achieved this goal. 3) Drawing upon the above reviews and discussions, we summarize ten prospective directions that present opportunities for deeper investigation into universal intelligence in healthcare.

There are some existing surveys related to the topic of multimodal learning in healthcare. For example, Shrestha *et al.* [8] explored the various aspects and advancements of medical vision language pre-training. Pei *et al.* [9] surveyed multi-modal learning on biomolecules, which primarily explores the technical advancements in integrating biomolecule sequences, 2D graphs, and 3D structures with natural language processing (NLP) techniques. Messina *et al.* [10] reviewed the RG for medical images and Zhao *et al.* [20] made a survey about relevant studies based on CLIP [17] architecture for medical imaging. Acosta *et al.* [7] surveyed multimodal biomedical AI, which mainly focuses

on the perspectives of data and applications while there are no technical details involved. However, it should be noted that these studies only summarize a specific portion of the entire multimodal research studies in healthcare and do not carry out in-depth discussions to explore whether current multimodal learning techniques can realize universal intelligent healthcare. To the best of our knowledge, this is the first comprehensive and systematic review on AI healthcare multimodal learning, covering a broad range of topics including datasets, task-oriented methodologies, contrastive FMs, multimodal large language models (MLLMs), as well as discussions on future directions. Through this panoramic and in-depth analysis, we conclude that current technology is not yet capable of achieving universal medical AI. We encourage researchers from different domains to collaborate and advance this important journey forward.

The following parts of the paper are organized as follows: (§2) gives preliminaries, including modalities, applications, and featured databases. (§3) shows the details of task-oriented methodologies for medical applications. (§4) and (§5) introduce details of two types of multimodal FMs, i.e., contrastive FMs and MLLMs, respectively. (§6) gives the discussions on five sub-questions to answer the proposed research question and (§7) further highlights ten future research directions. Finally, (§8) concludes the study.

2 PRELIMINARIES

2.1 Medical Modalities

Academically, modality refers to the way things are expressed or perceived, with every form or source of information being categorized as a modality [12]. In the healthcare domain, the term *multimodal* data typically pertains to digital records derived from diverse sources such as various machines, sensors, or experts, often represented in distinct formats. The medical modalities encompass elements such as medical vision, text, audio, and physiological signals, among others [7]. As shown in Figure 1, the medical vision modality consists of images obtained by different sensors, which are utilized for viewing the conditions or diseases

TABLE 1: Task formalization. ITR: image-to-text retrieval, TIR: text-to-image retrieval. \mathcal{R} refers to report and \mathcal{A} refers to answer. VQA_D/VQA_G denote discriminative/generative setting for VQA. TIP_{Seg} is a segmentation example and \mathcal{S} is the segmentation output. θ is the model parameter. \mathcal{Q} is the option set. \mathcal{Q} is the question and Γ is the database. $\tilde{\mathcal{V}}$ is the image of the other modality. \mathbf{G} . denotes whether the task is discriminative (◦) or generative (•).

| Tasks | \mathbf{G} . Formalization |
|-------------|---|
| RG | • $\mathcal{R} = [\hat{r}_1, \dots, \hat{r}_n]$, $\hat{r}_i = \arg \max_{w_i \in \mathcal{W}} \prod_{j=1}^i p(w_j \hat{r}_1, \dots, \hat{r}_{j-1}, \mathcal{V}; \theta)$ |
| VQA_D | ◦ $\mathcal{A} = \arg \max_{a_i \in \mathcal{O}} p(a_i \mathcal{V}, \mathcal{Q}; \theta)$ |
| VQA_G | • $\mathcal{A} = [\hat{a}_1, \dots, \hat{a}_n]$, $\hat{a}_i = \arg \max_{w_i \in \mathcal{W}} \prod_{j=1}^i p(w_j \hat{a}_1, \dots, \hat{a}_{j-1}, \mathcal{V}, \mathcal{Q}; \theta)$ |
| ITR | ◦ $\mathcal{T} = \arg \max_{t_i \in \Gamma_t} p(t_i \mathcal{V}; \theta)$ |
| TIR | ◦ $\mathcal{V} = \arg \max_{v_i \in \Gamma_v} p(v_i \mathcal{T}; \theta)$ |
| TIP_{Seg} | ◦ $\mathcal{S} = [\hat{s}_1, \dots, \hat{s}_k]$, $\hat{s}_i = \arg \max_{s_i \in \mathcal{O}} p(s_i \mathcal{V}, \mathcal{T}; \theta)$ |
| CIG | • $\mathcal{V} = \arg \max_{v_i} p(v_i \tilde{\mathcal{V}} / \mathcal{T}; \theta)$ |

of different organs or tissues. Within this scope, three types of images are commonly used, namely radiology, pathology, and camera images. Radiology is frequently employed to capture images of the human body’s internal conditions [21], primarily encompassing components such as X-ray, computed tomography (CT, 2D/3D), magnetic resonance imaging (MRI), positron emission tomography (PET), and ultrasound. Pathology is the scientific exploration of disease-induced alterations in cellular and tissue structures, which is conducted through the application of microscopy and supplementary laboratory methodologies [22], [23]. Beyond these image types, camera images provide a more direct depiction of the patient’s condition and are easier to gather, which is particularly effective and valuable in the detection of skin diseases [24]. Regarding medical text modalities, they generally encompass domain knowledge and data that are easy for humans to understand, collected from sources such as professional books, diagnostic reports, and literature. Recently, an increasing number of studies have focused on vision-language learning to provide complementary information and enhance performance [17].

Other medical modalities, such as audio, physiological signals (including electrocardiogram, *i.e.*, ECG, and electroencephalogram, *i.e.*, EEG), and electronic health record (EHR), also play significant roles in intelligent healthcare. Nonetheless, relevant studies predominantly concentrate on modeling the intricacies within these individual modalities, overlooking the interaction among multiple modalities. Thus, our work primarily centers on investigating multimodal studies involving medical images and text, and we discuss more comprehensive multimodal applications in §7.

2.2 Mainstream Applications

Formally, a medical image can be $\mathcal{V} \in \mathbb{R}^{c \times k \times h \times w}$, where c , k , h , w represent the channel, depth, height, and width of images, respectively. $k = 1$ denotes it is a 2D image and $k > 1$ indicates it is a 3D image. Language is represented as $\mathcal{T} = \{w^1, w^2, \dots, w^m\}$ with max token sequence m and w_i is the word token that from vocabulary \mathcal{W} . There are

various multimodal tasks for healthcare using intelligent technologies. For example, medical image fusion incorporates image features of different modalities. RG, VQA, cross-modal retrieval (image-to-text and text-to-image), text-augmented image processing (TIP), and cross-modal image generation (CIG) are common tasks in this cross-modal field. We summarize their formalization in Table 1 and illustrate them in Figure 2. Beyond these multimodal tasks, medical multimodal learning can also benefit unimodal tasks, such as medical image classification (IC), semantic segmentation (SS), and object detection (OD). By integrating data from various sources, medical multimodal learning improves feature representation, enhances contextual understanding, and supports data augmentation, thereby boosting the performance of these unimodal tasks. We will discuss studies for these applications in §4.

2.3 Featured Databases

There are outstanding medical recording databases, which are usually utilized for multimodal healthcare, such as PubMed³, MIMIC-CXR [25], and UMLS [26]. PubMed is a free medical literature database, gathering biomedical literature that composes medical journal articles, conference papers, and book chapters. It provides citations and abstracts, which may include links to the full text. Similarly, PubMed Central (PMC) provides an archive of full-text articles. MIMIC-CXR is a comprehensive database comprised of 377K images that correspond to a total of 227K chest radiographic studies. Each of these studies is accompanied by a detailed radiology report and pertinent chest X-ray (CXR) images. Authored by radiologists, these reports present a synopsis of their discoveries and typically consist of various sections, including examination, indication, impression, findings, technique, and comparison. UMLS, short for the unified medical language system, is a comprehensive collection of medical concepts from various lexicon resources. Each concept is assigned a unique identifier, which comes with corresponding definitions and numerous synonymous names. The UMLS also offers insights into the relationships between medical entities with a triplet format, conceptualizing a widespread medical knowledge graph (KG).

Using available resources, certain datasets are curated for specific tasks, such as the prevalent RG and VQA. We list some representative ones in Table 2 and Table 3, which are detailed described in the Appendix and further discussed in §6. They can be utilized for the data construction for MLLMs, as shown in §5. While we only discuss these two kinds of datasets, their utility extends to numerous tasks through the process of transformation. For instance, datasets used for RG can also be utilized for cross-modal retrieval, given their one-to-one correspondence.

3 MULTIMODAL MEDICAL STUDIES

3.1 Multi-modality Image Fusion

Images from a single modality provide limited insight into pathogenetic information within the human body. A reasonable fusion of multimodal medical images significantly

3. <https://pubmed.ncbi.nlm.nih.gov/>

TABLE 2: Representative RG datasets. # means the number of samples. *Col.* means the collection methods, where \circ , \star , and \bullet represent using synthetic, semi-automatic, and manual manner, respectively.

| Dataset | Time | #RG. | Images | | Source | Col. | Image & Report Characteristics |
|-------------------------|------|--------|--------|-----------|--------------------|-----------|---|
| | | | #Img. | Modality | | | |
| IU X-ray [27] | 2016 | 3.9K | 7.4K | CXR | In-House | \circ | <i>Indications, findings, impression, manual encoding, and MTI encoding.</i> |
| ICLEF-Caption-2017 [28] | 2017 | 184.6K | 184.6K | Multiple | PMC | \circ | Image captions in scholarly biomedical articles. |
| ICLEF-Caption-2018 [29] | 2018 | 232.3K | 232.3K | Multiple | PMC | \circ | Image captions in scholarly biomedical articles. |
| PEIR Gross [30] | 2018 | 7.4K | 7.4K | Pathology | PEIR | \circ | Images of gross lesions from sub-categories and one-sentence reports. |
| ROCO [31] | 2018 | 81.8K | 81.8K | Radiology | PMC | \circ | Reports are with UMLS CUIs/semantic types for image interrelations. |
| PadChest [32] | 2020 | 109.9K | 160.8K | CXR | In-House | \star | Reports contain findings, diagnoses, and locations in UMLS taxonomy. |
| MediCaT [33] | 2020 | 217K | 217K | Multiple | PubMed | \star | Containing captions, subfigures/subcaptions, and inline references. |
| ARCH [34] | 2021 | 11.8K | 15.1K | Pathology | PubMed & textbooks | \circ | Multiple instance captioning (a caption can relate to multiple images), including diagnostic, detection & classification, descriptive, etc. |
| FFA-IR [35] | 2021 | 10.7K | 1.0M | FFA | In-House | \bullet | Including bilingual (Ch. & En.) reports and explainable annotations. |
| CTRG [36] | 2024 | 2.8K | 8.1K | CT | In-House | \star | Reports (template/abnormal contents) of brain and chest CT scans. |

TABLE 3: Representative VQA datasets. # means the number of samples. *Col.* means the collection methods, where \circ , \star and \bullet represent using synthetic, semi-automatic and manual manner, respectively. *Gen.* denotes if the dataset is generative.

| Dataset | Time | #QA. | Images | | Source | Col. | Gen. | Characteristics & Contents |
|------------------------|------|--------|--------|-----------|------------------|-----------|--------------|--|
| | | | #Img. | Modality | | | | |
| VQA-Med-2018 [37] | 2018 | 6.4K | 2.8K | Radiology | PMC | \star | \checkmark | Rule-based question (location, finding, etc.) generation & experts check. |
| VQA-RAD [38] | 2018 | 3.5K | 315 | Radiology | MedPix | \bullet | \checkmark | About modality, plane, organ system, abnormality, object presence, etc. |
| VQA-Med-2019 [39] | 2019 | 15.2K | 4.2K | Radiology | MedPix | \circ | \checkmark | Test set is manually validated; about modality, plane, organ, abnormality. |
| VQA-Med-2020 [40] | 2020 | 5.0K | 5.0K | Radiology | MedPix | \circ | \checkmark | Test set is manually validated; about abnormality. |
| RadVisDial-Silver [41] | 2020 | 455.3K | 91.0K | CXR | MIMIC-CXR | \circ | \times | Four-choice questions; about 13 abnormalities. |
| RadVisDial-Gold [41] | 2020 | 500 | 100 | CXR | MIMIC-CXR | \bullet | \times | Four-choice questions generated by two radiologists. |
| PathVQA [42] | 2020 | 32.8K | 5.0K | Pathology | Textbooks & PEIR | \star | \checkmark | First dataset for pathology VQA, using a semi-automated pipeline; about color, location, appearance, shape, etc. |
| VQA-Med-2021 [43] | 2021 | 5.5K | 5.5K | Radiology | MedPix | \circ | \checkmark | Test set is manually validated; about abnormality. |
| SLAKE [44] | 2021 | 14.0K | 642 | Radiology | 3 datasets | \star | \times | Semantically annotated, knowledge-enhanced, and bilingual; about plane, modality, position, organ, KG, abnormal, shape, etc. |
| MIMIC-Diff-VQA [45] | 2023 | 700.7K | 164.3K | CXR | MIMIC-CXR | \circ | \checkmark | About abnormality, presence, view, location, type, level, and difference. |

contributes to a comprehensive understanding of intricate medical conditions [46], allowing clinicians to better delineate anatomical structures, lesions, and abnormalities. To establish a comprehensive view of how multimodal medical images are combined and analyzed, we introduce existing fusion strategies at the pixel, feature, and decision levels.

- **Pixel-level** fusion is a low-level fusion operation that concatenates pixels directly on the original image layer or their corresponding multi-resolution coefficients [47]. These approaches can be classified into three main categories: 1) multi-scale decomposition-based techniques [48]; 2) sparse representation methods [49], and 3) component substitution techniques [50]. However, the outcomes are often impacted by blurring effects that directly affect image contrast [51]. Moreover, there is a high registration requirement for multimodal images and it is usually sensitive to noise, making the pixel-level fusion process time-consuming and challenging.

- **Feature-level** fusion is a middle-level fusion strategy, and it is also recognized as the most commonly used strategy in deep learning methods. The common feature-level fusion strategy is to learn a shared representation or a joint embedding space from multiple features, using technologies such as adversarial learning [52], co-training [53], multi-task learning [54]. More recently, the transformer-based architectures, such as vision Transformer (ViT) [55], also show great versatility in handling different types of data, especially for heterogeneous data, and can be leveraged for feature-level fusion tasks by learning a joint representation. While feature-level methods overcome the drawbacks of pixel-based algorithms in terms of contrast, sensitivity to noise, and misregistration, they still have limitations such as spatial distortions [56]. For instance, medical images often

contain unclear regions due to poor illumination.

- **Decision-level** fusion is a high-level fusion strategy. It integrates multiple decisions derived from preliminary classifications and aggregates those decisions finally. Approaches are classified into two main categories: 1) hard fusion methods, which merge logical information membership values, such as model ensembling with majority or average voting [57]; and 2) soft fusion methods, where classifiers assign numerical values to reflect their confidence in decisions, or fuzzy classifiers are applied, such as fuzzy voting [58].

3.2 Medical Report Generation

Recently, researchers have proposed advanced strategies to enhance the quality and accuracy of automated medical RG, which can be categorized into three approaches: enhancing cross-modal alignment, improving through reinforcement learning techniques, and integrating auxiliary knowledge.

- **Cross-modal Alignment.** These studies focus on enhancing cross-modal alignment between medical images and reports to improve medical RG. Najdenkoska *et al.* [59] explored learning key topics between images and reports using variational topic inference to enhance semantic coherence. To facilitate multi-level cross-modal alignments, Li *et al.* [60] unified vision and text modalities into discrete tokens, which are then used to learn global semantic alignment and token-level alignment. Additionally, contrastive learning techniques [61] are also utilized for refined alignment between visual and textual data, enhancing the overall performance. For example, Wang *et al.* [62] introduced a phenotype-based contrastive learning framework that learns fine-grained representations, effectively bridging the gap between visual and textual modalities.

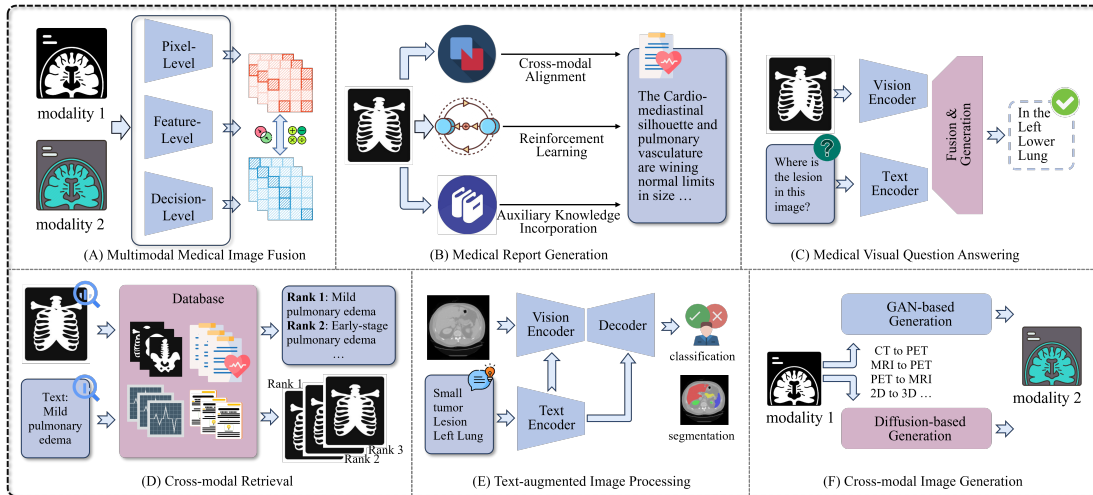


Fig. 2: Illustration of six mainstream types of task-oriented methodologies for medical applications.

- **Reinforcement Learning.** Reinforcement learning (RL) has gained significant traction in text generation due to its ability to use evaluation metrics as rewards, and updates model parameters via policy gradients. Inspired by this, researchers have applied RL to enhance medical RG models dynamically [63], using NLG (natural language generation) metrics-based rewards, *e.g.*, METEOR, ROUGE-L, BLEU, and CIDER. To improve factual correctness, Delbrouck *et al.* [64] employed RadGraph [65] to create semantic-based rewards for evaluating the factual accuracy and completeness of reports. Additionally, Parrer *et al.* [66] combined RL with text augmentation to enhance the quality and diversity of radiology reports, using BERTScore [67] and RadGraph [65] as rewards.

- **Auxiliary Knowledge Incorporation.** Recent research has utilized auxiliary signals like medical tags [68] and KGs to improve the coherence and accuracy of medical RG. For example, Li *et al.* [69] constructed dynamic KG to integrate both general and specific knowledge and introduced dynamic graph-enhanced contrastive learning to improve visual and textual representations. Huang *et al.* [70] developed the KiUT model, which integrates clinical knowledge by a symptom graph. It is combined with visual and contextual information through a U-Transformer architecture, significantly enhancing radiology RG. Hou *et al.* [71] proposed ORGAN, an observation-guided method using tree reasoning over observation graphs to improve interpretability and coherence. Additionally, researchers also included patients' prior examination data as anatomically aligned inputs to compare current and previous scans effectively [72], [73].

3.3 Medical VQA

The study of medical VQA focuses on interpreting medical images alongside spoken-language queries to generate accurate natural-language responses. It exhibits significant promise for diagnostic assistance and enhancing patient comprehension through educational support. The key research question revolves around the precise identification and comprehension of important regions within medical images (such as lesions, anomalies, and occupying masses),

and the semantic feature space alignment of these regions with the core demands expressed in the textual queries [74].

- **Knowledge Extraction Frameworks.** For encoder structures on a single modality, employing learning strategies that allow them to concurrently preserve both pre-existing general visual common sense and domain-specific medical knowledge represents an efficient model-agnostic solution. The first category employs the meta-learning paradigm, where the VQA model distributes the parameter training process across multiple related tasks to fuse loss gradients, ultimately adapting the model to medical VQA scenarios, *e.g.*, MAML [75] for model transfer and MMQ [76] for data refinement. Another approach to integrating medical knowledge with common sense is to perform conditional reasoning, which is designed to learn task-adaptive reasoning skills for different types of VQA tasks. This reasoning ability can be implemented by a supervised label loss on the embedding fusion of visual and textual encoders [77].

- **Pre-training & Fine-tuning Frameworks.** The remarkable advancements in pre-trained language models (LMs), such as BERT [78], have also introduced a novel paradigm for medical VQA. By designing self-supervised learning tasks, LMs can learn from vast text corpora without relying on external annotations, thus significantly reducing the cost of labeling. When faced with specific downstream VQA tasks, only a small amount of labeled data is required to finetune the model. In particular, knowledge relationships can also be injected into the pre-training process by non-euclidean encoders such as GNN [79]. Owing to the computational constraints of adjusting a large number of parameters during the fine-tuning phase, recent research has introduced VQA downstream task adaptation strategies that involve freezing parameters of the pre-trained model. For example, Liu *et al.* [80] proposed to establish a VQA adapter that is external to the pre-trained model, thus creating a plug-in for model adaptation to downstream tasks.

3.4 Cross-modal Retrieval

The medical cross-modal retrieval studies mainly line in the two categories: cross-modal retrieval within images, and

retrieval between images and texts. Cross-modal medical image retrieval involves searching for medical images in a database that have similar visual features to a given query image, thereby facilitating efficient clinical decision-making. Traditional studies handle the retrieval task by calculating similarities between texture features of different images, such as Radon Transform [81] on X-rays. Further, Mbilinyi *et al.* [82] suggested the application of deep features for extracting similar medical images from multimodal medical image databases. The results show that the retrieval performance of deep features obtained by CNNs is superior to the conventional texture features. Xu *et al.* [83] proposed multi-manifold deep discriminative cross-modal hashing for extensive medical image retrieval. The core aspect is the multi-modal manifold similarity that integrates multiple sub-manifolds based on heterogeneous data, thereby preserving the correlation among instances. This approach is both effective and efficient in adaptively retrieving medical images across various modalities. In general, this class of methods has undergone a transition from traditional feature extraction to efficient retrieval by deep neural networks.

Currently, the cross-modal retrieval between medical images and text, including ITR and TIR, is mainly done by learning embedded representations and calculating similarities by neural networks. The current major trend is to enhance representations with additional prior domain knowledge, such as the category information [84] and hierarchical semantic associations between disease labels [85].

3.5 Text-augmented Image Processing

Considering the complementary information contained in the text and image modalities, some studies focus on using text-guided information for medical image processing, including image classification and semantic segmentation.

Some text descriptions about quantity and scale could provide additional supporting information for image understanding and segmentation. For example, LViT [86] composes a U-shaped CNN branch and a U-shaped ViT branch for segmentation. It utilizes medical text annotation to address the limitations in image data quality, guiding the generation of enhanced pseudo labels in semi-supervised learning. Zhao *et al.* [87] and Dong *et al.* [88] introduced text-guided diffusion models for medical image segmentation. These models employ a text-attention mechanism to mitigate the influence of variations in the size and quantity of objects, such as representations of *one*, *many*, *small*, *medium*, and *large*, on the segmentation results. By concentrating on appropriate textual descriptors, the network is capable of adaptively modulating its focus towards the key characteristics of target objects, ensuring that the segmentation is both accurate and robust to changes in object attributes. Recently, medical FMs are frequently employed for text-guided image tasks, which is shown in §4 and §5.

3.6 Cross-modal Image Generation

In the medical domain, cross-modal image generation (also called modality translation) can be utilized in various scenarios, including education, data augmentation, missing data filling, and pathology research & understanding. From the technical perspective, they are categorized into two

classes: GAN-based (generative adversarial network) and diffusion-based. GAN-based models adopt the principle of adversarial training [89], involving two networks. The first generator is responsible for creating synthetic instances based on training data. The second discriminator is to differentiate between synthetic and real data. This competitive dynamic prompts the generator to produce highly realistic samples. In reality, CT scans emit radiation, potentially causing patient side effects, and their effectiveness in providing detailed images of soft tissue injuries is somewhat restricted. In contrast, MRI is radiation-free and safer. Thus, there is growing interest in generating CT images from corresponding MRI ones using GAN models, including perspectives of context-aware [90] and gradient consistency [91]. Also, there are several studies on the generation of other modalities, *e.g.*, CT to PET [92], MRI to PET [93], and PET to MRI [94].

Cross-modal diffusion-based generation models essentially transform the task of direct target generation into predicting random noise at every diffusion step. At its core, the diffusion model contains two critical processes: the forward diffusion process and the reverse denoising process. The forward diffusion process incrementally introduces Gaussian noise into an instance until it morphs into a sample of random noise. Conversely, the reverse denoising process aims to predict and remove the introduced noise [95]. Lyu *et al.* [96] made full use of denoising diffusion and score-matching strategies based on four different sampling approaches, implementing MRI to CT image synthesis. The results show that the model generates better synthetic CT images than the CNN and GAN models. Similarly, Meng *et al.* [97] introduced a unified multi-modal conditional score-based generative model to synthesize the missing modality using remaining modalities as conditions. The model employs only a score-based network to learn different cross-modal conditional distributions and the results show it can more reliably synthesize missing modality images of MRI.

4 CONTRASTIVE FOUNDATION MODELS (CFMs)

Given the intrinsic rarity, specificity, and specialized nature of data in the medical field, it is challenging and unrealistic to take large-scale high-quality annotated data for training. Therefore, some self-supervised strategies are introduced for building universal FMs [117]. FMs typically denote models that acquire the broad representation of general knowledge by undergoing pre-training on large-scale data through self-supervised learning. Subsequently, they can be refined through fine-tuning. Figure 3 illustrates the general architecture and applications of FMs in the medical domain. FMs have several key features: 1) pre-training on large generic datasets; 2) self-supervised learning strategies, such as contrastive learning and mask language modeling; 3) universal knowledge representation, meaning FMs learn a generic, task-independent knowledge representation that can be applied to a variety of different downstream tasks with a small amount of fine-tuning. According to the training strategies and applications, they can be categorized into two types: contrastive FMs (CFMs) and MLLMs. CFMs focus on learning a common cross-modal representation space by jointly optimizing the image encoder and text encoder to

TABLE 4: Representative CFMs. Abbreviations are as follows: *CPA*: cross-modal prototype alignment, *DEP*: disease/entity prediction, *ITM*: image-text Matching, *CTR*: cross-lingual text alignment regularization. “/” in the *Objectives & Training Process* splits different pre-training stages. Icons \star , $\color{red}\bullet$, and $\color{green}\bullet$ denote the module is frozen, updating, and inexistence when training, respectively. Their positions correspond image/adapter/language models. RN is short for ResNet.

| Model | Time | Modality | Image/Adapter/Language Model | Objectives (Training Process), Features & Applications |
|------------------|---------|-----------|-----------------------------------|---|
| ConVIRT [98] | 10/2020 | Radiology | RN50/-/ClinicalBERT | GCL ($\color{red}\bullet\color{red}\bullet\color{red}\bullet$); contrastive visual representation from paired descriptive text; IC, zero-shot ITR/TIR. |
| PubMedCLIP [99] | 12/2021 | Multiple | ViT-B-32, RN50/-/BioClinicalBERT | GCL ($\color{red}\bullet\color{red}\bullet$); fine-tuning CLIP on PubMed articles, followed by VQA-aware model; medical VQA. |
| CheXzero [100] | 09/2022 | CXR | ViT-B-32/-/Transformer | GCL ($\color{red}\bullet\color{red}\bullet$); self-supervised learning on CXR images; (low-resource) IC, <i>i.e.</i> , pathology detection. |
| BiomedCLIP [101] | 03/2023 | Multiple | ViT-B-16/-/PubMedBERT | GCL ($\color{red}\bullet\color{red}\bullet$); tuning on very large-scale dataset PMC-15M; (zero-shot) IC, ITR, TIR, VQA. |
| PLIP [22] | 03/2023 | Pathology | ViT-B-32/-/Transformer | GCL ($\color{red}\bullet\color{red}\bullet$); pathology data from Twitter; (zero-shot) IC, TIR, ITR, image representations. |
| PathCLIP [102] | 05/2023 | Pathology | ViT-B-16/-/Transformer | GCL ($\color{red}\bullet\color{red}\bullet$); tuning on PathCap with 207K high-quality image-text pairs; zero-shot IC, TIR. |
| CT-CLIP [103] | 03/2024 | 3D CT | CT-ViT/-/CXR-BERT | GCL ($\color{red}\bullet\color{red}\bullet$); 3D image abilities; (zero-shot) IC, volume-to-volume/report-to-volume retrieval. |
| PairAug [104] | 04/2024 | CXR | ViT-B-32/-/Transformer | GCL ($\color{red}\bullet\color{red}\bullet$); inter & intra data augmentation by ChatGPT and cross-attention maps; (zero-shot) IC. |
| GLoRIA [105] | 10/2021 | CXR | RN50/-/BioClinicalBERT | GCL+LCL ($\color{red}\bullet\color{red}\bullet$); text-attention local image representations; zero-shot ITR/IC, low-resource SS. |
| BioViL [106] | 04/2022 | CXR | RN50/-/CXR-BERT | GCL+MLM ($\color{red}\bullet\color{red}\bullet/\color{red}\bullet\color{red}\bullet$); language & vision-language tuning; (zero-shot) IC, SS, phrase grounding. |
| MedCLIP [107] | 10/2022 | CXR | ViT/-/BioClinicalBERT | soft GCL ($\color{red}\bullet\color{red}\bullet$); soft semantic matching for the false negative issue; (zero-shot) IC, ITR. |
| MGCA [108] | 10/2022 | CXR | ViT-B-16/-/BioClinicalBERT | GCL+LCL+CPA ($\color{red}\bullet\color{red}\bullet$); pathological region-level, instance-level & disease-level CL; IC, SS, OD. |
| BioViL-T [72] | 01/2023 | CXR | CNN-Transformer/-/CXR-BERT | GCL+LCL+MLM ($\color{red}\bullet\color{red}\bullet/\color{red}\bullet\color{red}\bullet$); multi-granularity and temporal connectivity modeling of images; |
| MedKlip [109] | 01/2023 | CXR | RN50/Transformer/ClinicalBERT | phrase grounding, (zero-shot, temporal) IC, RG, temporal sentence similarity. |
| KAD [21] | 02/2023 | CXR | RN50/Transformer/PubMedBERT | GCL+DEP ($\color{red}\bullet\color{red}\star$); incorporate entities & their descriptions; (zero-shot) IC, SS, region grounding. |
| PTUnifier [110] | 02/2023 | Radiology | CLIP-ViT-B/Transformer/RoBERTa | GCL/GCL+DEP ($\color{red}\bullet\color{red}\bullet/\color{red}\bullet\color{red}\star$); medical KG-enhanced; (zero-shot) IC, <i>i.e.</i> , disease prediction. |
| Med-UniC [111] | 05/2023 | CXR | RN50, ViT-B-16(L-32)/MLP/CXR-BERT | GCL+MLM+ITM ($\color{red}\bullet\color{red}\bullet$); soft prompts to unify early-fusion and later-fusion; ITR, TIR, VQA, etc. |
| MCR [112] | 12/2023 | CXR | ViT-B-16/-/BioClinicalBERT | GCL+CTR ($\color{red}\bullet\color{red}\bullet$); cross-lingual text alignment regularization for unifying cross-lingual (English & Spanish) medical vision-language pre-training; (zero-shot) IC, SS, OD. |
| MLIP [113] | 02/2024 | CXR | RN-50, ViT-B-16/-/ | GCL+MLM+MIM ($\color{red}\bullet\color{red}\bullet$); masked image and text as inputs, additional MLM and MIM; ITR, TIR. |
| MAVL [114] | 03/2024 | CXR | RN50/Transformer/ClinicalBERT | GCL+LCL+CPA ($\color{red}\bullet\color{red}\bullet$); multi-granularity, KG-based LCL; (zero-shot) IC, SS, OD. |
| KEP [115] | 04/2024 | Pathology | ViT-B-32(B-16)/-/PubMedBERT | GCL+DEP ($\color{red}\bullet\color{red}\star$); disease entities & their descriptions; (zero-shot) IC, SS, visual grounding. |
| DeViDe [116] | 04/2024 | CXR | ViT-B/LP/Med-KEBERT | AdaSP/GCL ($\color{red}\star\color{red}\bullet/\color{red}\bullet\color{red}\bullet$); incorporate a knowledge tree of disease attributes; ITR, TIR, disease retrieval, zero-shot classification on pathology patches, and zero-shot tumor subtyping on WSIs. |
| | | | | GCL ($\color{red}\bullet\color{red}\bullet$); entity-aware descriptions, global and local CL; (zero-shot) IC, IC, SS. |

maximize the similarity score of the positive sample (image-text pair) and minimize the similarity score of the negative sample [20]. However, MLLMs focus more on modeling the intrinsic cross-modal relationships, implementing cross-modal computation, and are capable of generating text outputs [18]. In this section, we will introduce CFMs and MLLMs will be detailedly elaborated in §5.

4.1 Overview of CFMs

Borrowing the idea of self-supervised contrastive learning (CL) that achieved great success in the computer vision field, CLIP [17] aligns vision and language semantic representations by pre-training on large-scale image-text pairs, which has greatly promoted the development of visual semantic understanding. It has sparked interest in its potential applications in the medical field. The general idea is to use an image encoder (*e.g.*, pre-trained ResNet [118] or ViT [55]) and a language encoder (*e.g.*, RoBERTa [119], CXR-BERT [106], ClinicalBERT [120], or PubMedBERT [121]) to obtain their corresponding representations. Rarely, an adapter serves as the bridge between them. Subsequently, they are updated using semantic contrast. We classify these studies into two categories: *CLIP-based* and *CLIP-variant* pre-training. The former generally uses typical global CL loss for modeling as like the original CLIP, while the latter introduces new optimizing objectives for additional specific concerns. Some representative medical CFMs are listed in Table 4 and some representative datasets for alignment are in Table 6.

4.2 CLIP-based Pre-training

Formally, given the image set $\{\mathcal{V}_1, \mathcal{V}_2, \dots, \mathcal{V}_N\}$ and text set $\{\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_N\}$ where \mathcal{V}_i and \mathcal{T}_i form a image-text pair, their representations $\mathbf{v} \in \mathbb{R}^{N \times n \times d}$ and $\mathbf{t} \in \mathbb{R}^{N \times m \times d}$ can be obtained by image and text encoders, respectively. This formalization is similar to §2.2. Further, their global

representation can be obtained by pooling operation or taken a representative one as $\mathbf{v}^g \in \mathbb{R}^{N \times d}$ and $\mathbf{t}^g \in \mathbb{R}^{N \times d}$. \mathcal{D} , \mathcal{N} , and \mathcal{M} are the index sets for the samples, image patches, and text tokens. Based on them, global GL (GCL) employs a comprehensive and holistic perspective for semantic relationships. Its loss is to directly model the whole semantics between image and text using InfoNCE loss:

$$\mathcal{L}_{\text{GCL}} = \mathbb{E}_{i \in \mathcal{D}} \left[-\log \frac{\exp(s(\mathbf{v}_i^g, \mathbf{t}_i^g))/\tau}{\sum_{j=1}^N \exp(s(\mathbf{v}_i^g, \mathbf{t}_j^g))/\tau} \right]. \quad (1)$$

s denotes the cosine similarity function and τ is a pre-set temperature parameter. Note that it may $s(a, b) \neq s(b, a)$, so $s(\mathbf{t}_i^g, \mathbf{v}_i^g)$ may also be calculated for comprehensive modeling. Using GCL, many studies learn semantic representations for medical images and their corresponding description texts. For example, studies on CXR (CheXzero [100] & PairAug [104]), pathology (PLIP [22] & PathCLIP [102]), radiology (ConVIRT [98]), and multiple modalities (PubMedCLIP [99] & BiomedCLIP [101]). Unlike these approaches concentrating on 2D images, Hamamci *et al.* [103] proposed the first 3D medical imaging dataset CT-RATE with textual reports. Based on it, CT-CLIP is pre-trained using a 3D encoder for chest CT volume representations and it is then aligned with CXR-BERT outputs.

4.3 CLIP-variant Pre-training

Beyond GCL, there are studies with other modeling objectives, which can be summarized as following four types.

- **Intra-modal modeling.** Beyond inter-modal modeling, there are additional objectives focused on intra-modal modeling, where two standard objectives are typically employed. MIM (mask image modeling) [112] usually uses the mean squared error function to compute the normalized pixel-wise difference between the original image patches P_i and reconstructed image patches \tilde{P}_i . MLM (mask language

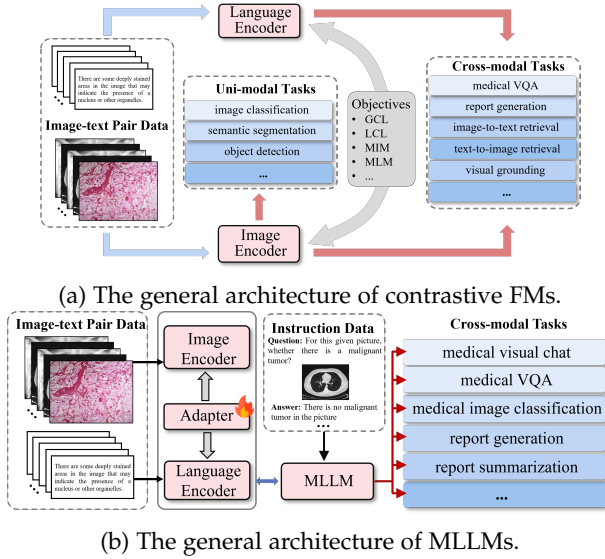


Fig. 3: Illustration of FMs in the medical domain.

modeling) [72], [110], [112] are to predict masked tokens based on other given tokens. They are calculated as follows:

$$\mathcal{L}_{\text{MIM}} = \mathbb{E}_{i \in \mathcal{D}} [(P_i - \tilde{P}_i)^2], \quad \mathcal{L}_{\text{MLM}} = \mathbb{E}_{i \in \mathcal{D}, j \in \mathcal{M}} [\mathbb{I}(\mathcal{T}_i^j) \cdot f(\mathcal{T}_i^j)]. \quad (2)$$

\mathbb{I} is to indicate where token T_i^j is masked, where $\mathbb{I}(\mathcal{T}_i^j) = 1$ if masked. Otherwise, the value is 0. f is the loss for the predicted token compared to the original one. These two approaches acquire internal semantics of modality through reconstruction, enhancing single-modal model capability.

- **Multi-granularity.** Focusing on the global information of images and texts may not be sufficient for capturing fine-grained information. So based on GCL, its variant of multi-granularity is introduced, where local CL (LCL) is the representative. It utilizes the global (or local) representations of one modality to align local representations of another. Local token representations of text and image are \mathbf{t}_i^j and \mathbf{v}_i^j , which means the token j of sample i . The corresponding global representation can be $\tilde{\mathbf{t}}_i^g$ and $\tilde{\mathbf{v}}_i^g$, obtained by pooling strategy or linear transformation based on \mathbf{t}_i^g and \mathbf{v}_i^g . There are two main LCL approaches (global-to-local as example):

$$\mathcal{L}_{\text{LCL(I2T)}} = \mathbb{E}_{i \in \mathcal{D}, j \in \mathcal{M}} \left[-\log \frac{\exp(\sigma(\tilde{\mathbf{v}}_i^g, \mathbf{t}_i^j)) / \tau}{\sum_{k=1}^N \exp(\sigma(\tilde{\mathbf{v}}_k^g, \mathbf{t}_i^j)) / \tau} \right], \quad (3)$$

$$\mathcal{L}_{\text{LCL(T2I)}} = \mathbb{E}_{i \in \mathcal{D}, j \in \mathcal{N}} \left[-\log \frac{\exp(\sigma(\tilde{\mathbf{t}}_i^g, \mathbf{v}_i^j)) / \tau}{\sum_{k=1}^N \exp(\sigma(\tilde{\mathbf{t}}_k^g, \mathbf{v}_i^j)) / \tau} \right]. \quad (4)$$

The former uses global image information as the anchor, aligning it with corresponding local text counterparts. MGCA [108] and BioViL-T [72] follow the latter manner, viewing the global text representation as the anchor. GLORIA [105] uses both manners, where local representations are obtained by text token-based attentive image pooling.

In a different manner, MLIP [113] introduces local token-knowledge-patch alignment using a medical KG, *i.e.*, UMLS. The cross-attention is used to match knowledge and image or text, with the input of pre-trained entity embeddings of UMLS and feature of image patch or text token. Processed representations are compared with the original ones for

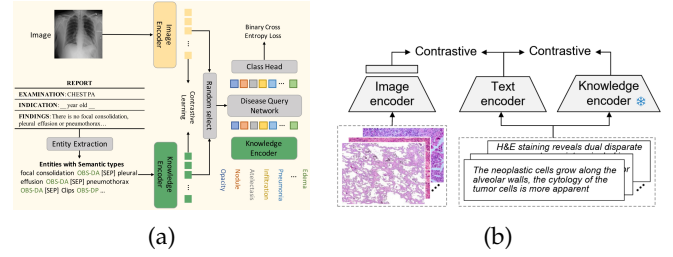


Fig. 4: Illustration of CFM pre-training with structural knowledge fusion. (a) The second pre-training stage of KAD utilizes contrastive loss to align image and entity content representations and employs cross entropy for disease prediction. Taken from [21]. (b) The second pre-training stage of KEP. Taken from [115].

knowledge alignment, from both image and text sides. Thus, MLIP is in a manner of local-local contrast rather than other global-local methods. Beyond LCL, MGCA [108] also introduces cross-modal prototype alignment, which assumes that images with the same disease would have similar disease-level representations. It has a higher level than the global instance and the local tokens. To realize it, MGCA pre-defines trainable embeddings for cross-modal prototypes with a certain number, which can be used for *pseudo-label* calculation with global text and image representations. Finally, the model is optimized by aligning the text and image *pseudo-labels*. Similarly, MLIP also introduces a higher level, *i.e.*, category-level, where KG embeddings and original features are guided for category clustering.

- **Structural Knowledge Fusion.** To inject professional knowledge of the medical domain to enhance the semantic representations, knowledge-fused pre-training is proposed. They usually have two pre-training stages, where the first is to incorporate domain-structured knowledge to optimize the knowledge encoder and the second is for image-text alignment with the pre-trained knowledge encoder. KAD [21] leverages medical knowledge to guide vision-language pre-training. A well-established medical KG, *i.e.*, UMLS, is introduced to fine-tune PubMedBERT by contrastive loss of concept-definition pairs and concept-relation-concept triplets. Then the given raw X-ray reports are converted into contents of medical entities and their presence, using heuristically defined rules, RadGraph [65], or ChatGPT. The pre-trained knowledge encoder is used to guide the visual representation learning by contrastive learning between representations of image and generated entity contents, effectively injecting the domain knowledge into the visual encoder. Besides, the query disease is also input to incorporate randomly selected image or text entity content representations for disease prediction. In the inference stage, by inputting unseen diseases, KAD can handle zero-shot disease prediction given an image. KEP [115] curates a pathology *knowledge tree* PathKT, consisting of three-level tree structures: tissue, disease, and attribute. Each disease entity corresponds to several attributes, including disease synonyms, definitions, cytology and pathology features. The knowledge encoder is updated by metric learning with AdaSP loss [122], such that the representations of a specific disease and its attributes are close in the embedding

space. After that, the text encoder (initialized with the weights of the knowledge encoder) and image encoder are updated by image-text contrastive learning. These methods of structural knowledge fusion are illustrated in Figure 4.

- **Other Variants.** For various aspects of multimodal medical applications, some other variants are proposed. Previous methods could encounter many false negatives, meaning that images and reports from different patients could have the same semantics but are mistakenly treated as negative samples. So MedCLIP [107] decouples image text pairs and conducts contrastive learning to reduce false negatives by introducing external medical knowledge. To make full use of limited usable data and fix false negatives in contrastive learning, MedCLIP introduces UMLS to detect 14 main entity types for images with diagnosis labels. Multi-hot vectors of 14 dimensions from the extracted entities are obtained for images and texts, from which the semantic similarity is calculated. The semantic similarity is viewed as *soft targets* for model training rather than 0/1 labels in the original contrastive loss. In this way, unpaired data can also be taken into consideration. To make the model capable of temporal information in the medical domain, BioViL-T [72] exploits *temporal correlations* by making prior images available for comparison to a given report. The visual representations of the two images combine to make global and local contrastive learning, where an additional MLM is utilized for text-side pre-training. PTUnifier [110] introduces the *soft prompts* to unify early-fusion and later-fusion medical vision-language pre-training, making it compatible with different kinds of inputs, including image-only, text-only, and image-text pairs. It constructs prompt pools for different modalities so that different inputs can select their corresponding prompts, improving the prompt diversity and the model scalability. To consider the presence of community bias caused by different languages, Med-UniC [111] unifying *cross-lingual* (English & Spanish) medical multi-modal by diminishing bias. Besides GCL for the vision-vision and vision-language alignment, cross-lingual text alignment regularization, including text augmentation, text-feature alignment, and text-to-text alignment, learns language-independent text representations and neutralizes the adverse effects of community bias on other modalities.

4.4 Data Augmentation

Medical texts are usually characterized by their specialized and condensed nature, making them difficult to understand by layman and neural models. Therefore, several studies have introduced augmented text descriptions. MedKLIP [109] focuses on the entities in the medical reports and adds entity descriptions. The representations of these added descriptions are fused with image features to make predictions of entity existence and its location. Based on it, MAVL [114] further expands the description of disease entities to multiple visual aspects, including pattern, texture, opacity, border, location, shape, and fluid presence, where GPT-4 is utilized to programmatically generate descriptions of these aspects. Beyond the loss functions of MedKLIP, it introduces another contrastive target to align visual representation and that of each entity aspect’s description, empowering MAVL with the ability of zero-shot recognition

of unseen diseases. DeViDe [116] utilizes publicly-available Mixtral-8x7B [123] to collect and process radiographic descriptions, for a specific entity or a disease.

Considering most augmentation techniques tend to narrow their focus, prioritizing either text or image augmentation, rather than blending the two. PairAug [104] designs a pairwise augmentation approach that contains an inter-patient augmentation (InterAug) branch and an intra-patient augmentation (IntraAug) branch. Specifically, the InterAug branch generates radiology images using synthesised yet plausible reports derived from an LLM. IntraAug branch uses newly generated reports to manipulate images. This process facilitates the generation of new paired data for each individual with diverse medical conditions, where ChatGPT is used to report modification.

4.5 Downstream Applications

Based on pre-trained CFMs, many medical applications can be achieved, varying from uni-modal to cross-modal tasks.

- **Uni-modal Tasks.** The jointly pre-trained image encoder and text encoder can be individually or jointly used for many uni-modal tasks, such as image classification, semantic segmentation, and object detection. Medical image classification is usually to detect diseases in the image, based on GCL pre-training, CFMs can do zero-shot image classification by calculating the similarity between image representation and text prompts with a specific disease, *e.g.*, *this is an image of {disease}* and *{disease} presented in the image* [101]. Also, the image encoder can be frozen and the subsequent MLP is updated to fine-tune for image classification tasks, namely linear probing. For semantic segmentation and object detection tasks, the pre-trained image encoder is initialized as the backbone encoder, followed by a trainable task-specific decoder, like U-Net [124] or ResUNet [125] and YOLOv3 [126] for these two tasks, respectively.

- **Cross-modal Tasks.** Pre-trained CFMs can also be used for many cross-modal tasks, including VQA, RG, ITR, TIR, and visual grounding. As the CFMs generally have no ability for text generation, so its application for VQA mainly focuses on the classification setting, using VQA-RAD and SLAKE datasets [99], [110]. The VQA model is usually initialized with a pre-trained CFM encoder and incorporates an MLP or specific decoder to make predictions. For report generation, BioViL-T processes the prior report and both images (prior and current) with an encoder and an additional decoder is utilized for generation, where two broad categories, *i.e.*, nearest-neighbour and auto-regressive can be used. Retrieval-based tasks (ITR, TIR and disease retrieval) are directly realized by calculating similarities between the two modalities’ representations. Similarly, phase grounding calculates the similarity of representations between the text phase and the local image patch [72], [106].

5 MULTIMODAL LLMs (MLLMs)

Benefiting from the rapid development of LLMs, MLLMs, also known as visual language models (VLMs), have garnered significant attention from researchers owing to their powerful representational capabilities and remarkable proficiency in handling multimodal data [18]. Their general

TABLE 5: Representative MLLMs in the general and medical domain. “/” in the *Datasets & Training Process* splits different pre-training stages. Icons \star , $\color{red}\star$, and $\color{green}\star$ denote the module is frozen, updating, and inexistence when training, respectively. Their positions correspond image/adapter/language models. RN is short for ResNet. † means to omit the pure LLM pre-training or the pure image encoder pre-training.

| Model | Time | Modality | Image/Adapter/Language Model | Size | Datasets (Training Process) & Contribution |
|--------------------|---------|-----------|------------------------------|-----------|---|
| Flamingo [127] | 04/2022 | Natural | NFNet/QR+CA/Chinchilla | 3/9/80B | M3W, ALIGN, LTIP, VTP ($\star\star\star$); interleaved visual-textual data, few-shot in-context learning abilities. |
| CoCa [128] | 05/2022 | Natural | ViT/CA/Transformer | 2.1B | JFT-3B, ALIGN ($\color{red}\star\star\star$); CCL for image/text encoder, additional text decoder for language generation. |
| BLIP-2 [129] | 01/2023 | Natural | ViT/QR/OPT, FlanT5 | 3.1-12.1B | COCO, CC3M, etc ($\star\color{green}\star/\star\color{red}\star$); QFormer with cross-modal tasks, two-stage bootstrapping strategy. |
| LLaVA [130] | 04/2023 | Natural | ViT-L-14/LP/Vicuna | 13B | CC3M/LLaVA-Instruct-158K, ScienceQA ($\star\color{green}\star/\star\color{red}\star$); instruction-following data using GPT-4. |
| MiniGPT-4 [131] | 04/2023 | Natural | ViT-G-14/LP+QR/Vicuna | 13B | Conceptual Caption, SBU, LAION, etc ($\star\color{red}\star/\star\color{red}\star$); only update a LP layer for alignment. |
| SkinGPT-4 [24] | 04/2023 | Camera | ViT-G-14/LP+QR/Vicuna | 13B | SKINCON, Dermnet ($\star\color{red}\star/\star\color{red}\star$); based on MiniGPT-4, skin disease diagnoses by uploading photos. |
| PathAsst [102] | 05/2023 | Pathology | ViT-B-16/LP/Vicuna | ~13B | PathInstruct ($\star\color{red}\star/\star\color{red}\star$); specific model for pathology, capable of invoking eight sub-models. |
| MedBLIP [132] | 05/2023 | 3D MRI | ViT-G-14/QR/BioMedLM | ~2.7B | ADNI, NACC, etc ($\color{red}\star\star\star/\star\color{red}\star$); based on BLIP-2, LoRA, 2D vision encoder for 3D MRI scans. |
| LLM-CXR [133] | 05/2023 | CXR | VQ-GAN (RN)/-/Dolly-v2-3B | ~3B | MIMIC-CXR ($\color{green}\star\color{green}\star/\star\color{green}\star$); no adapter, image tokens (input or output), CXR generation abilities. |
| BiomedGPT [134] | 05/2023 | Multiple | VQ-GAN (RN)/-/BART | 33-182M | 14 datasets ($\color{green}\star\color{green}\star/\star\color{red}\star$); no adapter, unified FM, zero-shot transfer learning, diverse biomedical tasks. |
| XrayGPT [135] | 06/2023 | CXR | MedCLIP/LP/Vicuna | 13B | MIMIC-CXR, OpenI ($\star\color{red}\star/\star\color{red}\star$); based on MiniGPT-4, answer open-ended questions about CXR. |
| LLaVA-Med [136] | 06/2023 | Multiple | ViT-L-14/LP/Vicuna | 7/13B | LLaVA-Med-Align/LLaVA-Med-Inst ($\star\color{red}\star/\star\color{red}\star$); based on LLaVA, instructions using GPT-4. |
| Med-Flamingo [137] | 07/2023 | Multiple | ViT-L-14/QR+CA/LLaMA | 8.3B | MTB, PMC-OA ($\star\color{red}\star$); based Flamingo, few-shot in-context learning abilities, multi-image input ability. |
| Med-PaLM M [138] | 07/2023 | Multiple | ViT/LP/PaLM | 12-62B | MultiMedBench ($\color{red}\star\star\star$); generalist biomedical AI system, closed-source, can handle multiple modalities. |
| RaDfM [139] | 08/2023 | Radiology | 3D ViT/QR/MedLLaMA | 14B | MedMD/RedMD ($\color{red}\star\star\star/\star\color{red}\star$); 2D&3D, tuning on very large-scale datasets, multi-image input ability. |
| RaDialog [140] | 10/2023 | CXR | RN50/QR/Vicuna | ~7B | MIMIC-CXR/image-grounded instruct data ($\star\color{red}\star/\star\color{red}\star$); LoRA, radiology RG & interactive dialog. |
| Qilin-Med-VL [141] | 10/2023 | Multiple | ViT-L-14/LP/Chinese-LLaMA2 | ~13B | ChiMed-VL-Align/ChiMed-VL-Inst ($\star\color{red}\star/\star\color{red}\star$); Chinese medical MLLM, multiple image modalities. |
| MAIRA-1 [142] | 11/2023 | CXR | RAD-DINO/LP/Vicuna-7B | ~7B | MIMIC-CXR ($\star\color{red}\star$); larger image resolution (518x518), leveraging GPT-3.5 for data augmentation. |
| PathChat [143] | 12/2023 | Pathology | ViT-L/LP+QR/LLaMA2 | ~13B | Alignment/instruction dataset ($\star\color{red}\star/\star\color{red}\star$); vision-language generalist assistant for human pathology. |
| MedXChat [144] | 12/2023 | CXR | ViT-L-14/-/LLaMA | ~7B | MIMIC-CXR ($\star\color{red}\star$); no adapter, LoRA, finetuned Stable Diffusion model for text-to-CXR synthesis. |
| CheXagent† [145] | 01/2024 | CXR | EVA-CLIP-g/LP+QR/Mistral | 8B | MIMIC-CXR, CheXInstruct, etc ($\color{green}\star\color{green}\star/\star\color{red}\star/\star\color{red}\star$); multi-stage training, systematical evaluation. |
| CONCH [23] | 03/2024 | Pathology | ViT-B-16/CA/Transformer | - | PubMed, Internal Data ($\color{red}\star\star\star$); based on CoCa, using diverse sources of histopathology images. |
| M3D-LaMed [146] | 03/2024 | 3D CT | 3D ViT/Pooling+LP/LLaMA2 | 6.9B | M3D-Data ($\star\color{red}\star$); a versatile 3D MLLM, extra segmentation module for direct 3D image segmentation. |
| Dia-LLaMA [147] | 03/2024 | 3D CT | ViT3D/QR/LLaMA2 | ~7B | CTRG-Chest-548K ($\star\color{red}\star$); LoRA, disease-aware attention, additional diagnostic information for RG. |
| LLaVA-Rad [148] | 03/2024 | CXR | BiomedCLIP/LP/Vicuna | 7B | CXR-1M, MIMIC-CXR ($\star\color{red}\star/\star\color{red}\star$); LoRA, multi-stage tuning for lightweight MLLM. |
| WoLF [149] | 03/2024 | CXR | CLIP-ViT-L-14/?/Vicuna | ~7B | MIMIC-IV, MIMIC-CXR ($\star\color{red}\star/\star\color{red}\star$); health-specific instruction and anatomy-specific knowledge. |

modeling objectives are the next token prediction based on the image and previous text tokens:

$$\mathcal{L}_{\text{MLLM}} = \mathbb{E}_{i \in \mathcal{D}, j \in \mathcal{M}} -\log p(\mathcal{T}_i^j | \mathcal{V}_i, \mathcal{T}_i^{< j}). \quad (5)$$

From both theoretical and application standpoints, there exist distinct differences between CFMs and MLLMs: 1) CFMs are generally tuned based on the image-text pair data, whereas MLLMs focus on the multimodal instruction-following data; 2) GCL is the main objective for the CFMs while MLLMs are to generate text based on multimodal inputs; 3) CFMs are typically used for discriminative tasks, but MLLMs are more commonly used for generative tasks. We summarize some typical general MLLMs and noted medical MLLMs in Table 5.

5.1 Modality Encoder and Cross-modal Adapter

MLLMs employ a straightforward yet effective method for building FMs. This approach involves the use of an image encoder and language model, collectively known as modality encoders, which facilitates corresponding representations in latent spaces. Additionally, a cross-modal adapter is introduced to align these representations of various modalities within a shared space.

• **Image Encoder.** Pre-trained image encoders are typically employed to generate image representations, which are subsequently integrated with LLMs for multimodal tasks. Commonly utilized pre-trained encoders include NFNet [150], ViT [55], CLIP ViT [17], EVA-CLIP [151] of the general domain. Besides, to enhance the encapsulation of medical knowledge within the latent representations, pre-trained medical image encoders are employed in the creation of medical MLLMs. For instance, the vision components of image-text pre-trained PathCap [102], BioMedCLIP [101], and BioViL-T [72] act as the image encoders of

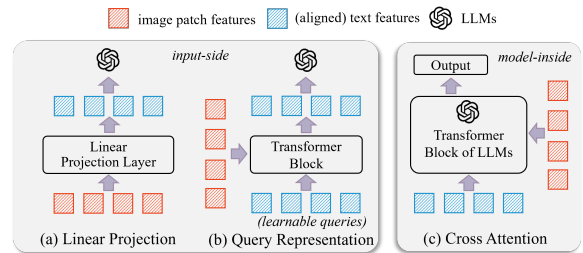


Fig. 5: Illustration of different types of adapters for MLLMs. LR and QR are input-side while CA is inside of LLMs.

PathAsst [102], LLaVA-Med [136], and RaDialog [140], respectively. Additionally, RAD-DINO [152] employs CXR image pre-training, based on the self-supervised pre-training strategy of the DINO model [153].

• **Language Model.** Building on the powerful LLMs within the NLP domain, MLLMs employ them to process textual input and generate corresponding textual responses. Both general-purpose and medical-specific MLLMs utilize the Transformer-style architecture as the text encoder. This is evident in models such as OPT [154], Flan-T5 [155], Vicuna [156], Mistral [157], LLaMA [158], LLaMA2 [159], and the Chinese-LLaMA2 [160], which is catered to the general domain. In addition, there are several models tailored to medical languages, such as BioGPT [161], BioMedLM [162], and MedLLaMA [163], also in widespread use.

• **Cross-modal Adapter.** The cross-modal adapter serves as a connector between image and text representations within MLLMs. Primarily, it encompasses three categories: *Linear Projection*, *Query Representation*, and *Cross Attention*, as illustrated in Figure 5. The first two types are handled during the input stages, transforming the hidden image representation into virtual token embeddings that align with text token embeddings. The last type typically manifests within the

internal computational procedures of the LLM [18].

Linear Projection (LP). Typically, it employs a one or two-layer MLP to transform the image embedding space into a textual one. Simple yet effective, LLaVA [130] first introduces it in the general MLLMs and has since been adopted by medical MLLMs, such as PathAsst [102], XrayGPT [135], LLaVA-Med [136], and MAIRA-1 [142].

Query Representation (QR). This approach establishes learnable query representations for images, which are then combined with textual token representations as inputs for LLMs. The quantity of queries is usually significantly less than the number of image patches, thereby reducing computational efforts and enhancing efficiency. The *perceiver resampler* in Flamingo [127] is designed to handle a flexible number of visual varying-size features (typically large) that are obtained from the vision encoder, generating a reduced number of visual outputs. Similarly, QFormer, *i.e.*, Querying Transformer, proposed in BLIP-2 [129] is to unify the image information in the language space. It first introduces several learnable queries that interact with image features using cross attention. To align the multimodal input, three pre-training objectives are usually used before fine-tuning, *i.e.*, image-text contrastive learning, image-grounded text generation, and image-text matching. Regardless of the size of the visual encoder, the QFormer’s final output length remains constant, *e.g.*, 32 for BLIP-2 and MedBLIP [132], significantly diminishing the computational load.

Cross Attention (CA). Inspired by the self-attention of Transformer, the cross attention approach assigns varying roles to image and language representations as *query*, *key*, or *value* during the computational process of the LLM. For example, the *GATED XATTN-DENSE* in Flamingo [127] views image representations as the *key* and *value*, the language counterpart as the *query*. Similarly, CoCa [128] and CONCH [23] employ cross attention for image information incorporation in the text decoder.

5.2 Tuning Process & Technical Details

Based on the image encoder and language models, the implementation of MLLMs is achieved through the use of *pre-training* and *instruction-tuning* strategies [8]. The primary objective of pre-training is typically to adapt initializations from the general domain to the specialized medical one and to bridge distinct modalities. Instruction refers to the task description and the goal of instruction tuning is to enhance a model’s comprehension of user instructions and execute the tasks accordingly. Through this process, MLLMs are capable of generalizing to previously unseen tasks with new instructions, thereby enhancing zero-shot performance [18]. As the representative study, LLaVA-Med [136] is first fine-tuned on LLaVA-Med-Align (600K biomedical image-text pairs) to update linear transformation layer and then carries out the instruction-tuning on LLaVA-Med-Inst (60K image-text responses collected using GPT-4). For parameter-efficient tuning, the LoRA [164] technique, *i.e.*, low-rank adaptation, could be utilized [132], [140], [147], [148].

General MLLMs. Beyond the contrastive loss in CLIP, CoCa [128] introduces the caption loss in the architecture. The model comprises an image encoder, a text decoder, and a multimodal text decoder, capable of handling text

generation tasks. Flamingo [127] is groundbreaking research in the realm of general MLLM, which has been trained on vast multimodal web corpora containing arbitrarily *interleaved* text and images. This approach is vital for equipping Flamingo with in-context few-shot learning abilities. In its framework, both the pre-trained vision encoder and language model are frozen. Two novel components, namely the *perceiver resampler* and *GATED XATTN-DENSE*, are incorporated to effectively bridge the gap between powerful vision-only and language-only models. The *GATED XATTN-DENSE* block is integrated into every layer of the frozen LLM, facilitating cross attention between vision and language features. BLIP-2 [129] pre-trains a lightweight QFormer following a two-stage strategy to bridge the modality gap. It first bootstraps vision-language representation learning from a frozen image encoder and the second stage bootstraps vision-to-language generative learning from a frozen LLM.

LLaVA [130] and MiniGPT-4 [131] are the pioneers of MLLMs tuning with image-text instruction data. LLaVA is an initial endeavour to employ a language-only GPT-4 model for the creation of multimodal language-image instruction data known as LLaVA-Instruct-158K. To effectively leverage the capabilities of both pre-trained LLM and visual model, a linear projection is introduced to align the vision features to the language counterpart. LLaVA is pre-trained with only 1 epoch for feature alignment, where only parameters of the linear projection are optimized. Then, keeping the visual encoder weights frozen, both pre-trained weights of the projection layer and LLM are updated on LLaVA-Instruct-158K for 3 epochs or on ScienceQA [165]. To realize numerous advanced multi-modal abilities demonstrated by GPT-4, MiniGPT-4 aligns a frozen visual encoder with the frozen Vicuna model using QFormer and a projection layer.

Medical MLLMs. Motivated by the remarkable achievements of general MLLMs, medical MLLMs are developed with the aim of creating a highly efficient universal medical assistant. For example, SkinGPT-4 [24], LLaVA-Med [136], and Med-Flamingo [137] are the adaptations of MiniGPT-4, LLaVA, and Flamingo within the medical field, respectively. For human pathology, PathAsst [102] and PathChat [143] are developed through instruction tuning, utilizing pathology image-text data. Taking PathAsst as an example, it constructs instruction-following data PathInstruct, which contains description-based and conversation-based, to revolutionize diagnostic and predictive analytics in pathology. ChatGPT is utilized to generate conversational QA pairs based on image captions. Special model-invoking instruction-following samples are also included. Based on two-step tuning, PathAsst has the ability for VQA and conversation tasks related to pathology. Also, it can invoke sub-models for more comprehensive applications, such as LBC (liquid-based cytology) classification and LBC detection.

Actually, a majority of medical MLLMs are predominantly focused on CXR or radiology modalities. This concentration arises from the wealth of data available, exemplified by databases such as MIMIC-CXR. A number of noteworthy studies exist in this direction, including LLM-CXR [133], XrayGPT [135], MAIRA-1 [142], MedXChat [144], CheXagent [145], LLaVA-Rad [148], and WoLF [149]. These are constructed in line with the general strategies of MLLMs.

TABLE 6: Representative datasets for medical FMs, including for alignment and instruction-tuning. Column *LLM* indicates whether using LLM for dataset construction.

| Dataset | Modality | Scale | LLM | Remarks | Source | Time |
|-----------------------|-----------|-------|-----|---|----------------------------|---------|
| PMC-15M [101] | Multiple | 15M | ✗ | collected from scientific papers, spanning thirty major biomedical image types. | PMC | 03/2023 |
| OpenPath [22] | Pathology | 208K | ✓ | employ 32 pathology-related hashtags to retrieve relevant tweets. | Twitter, PathLAION | 03/2023 |
| PathCap [102] | Pathology | 207K | ✓ | pathology image-caption pairs, ChatGPT is employed to refine the captions. | PubMed, Books, Cytologists | 05/2023 |
| LLaVA-Med-Align [136] | Multiple | 600K | ✗ | image-text pairs sampled from PMC-15M, images of multiple modalities for alignment. | PMC-15M | 06/2023 |
| MedMD [139] | Radiology | 16M | ✗ | multiple modalities of 2D&3D, covering a wide range of anatomies with over 5000 diseases. | 14 Sources | 08/2023 |
| ChiMed-VL-Align [141] | Multiple | 580K | ✓ | covering X-ray, MRI, CT, radioisotope, mitotic, etc, translated by GPT-3.5 to Chinese. | PMC-OA, PMC-CaseReport | 10/2023 |
| CT-RATE [103] | 3D CT | 50K | ✗ | 3D chest CT volumes and corresponding radiology text reports from hospital. | Internal Hospital | 03/2024 |
| PathInstruct [102] | Pathology | 180K | ✓ | description-based & conversation-based (via ChatGPT), containing model-invoking part. | PathCap, Human Design | 05/2023 |
| LLaVA-Med-Inst [136] | Multiple | 130K | ✓ | using GPT-4 to generate multi-round Q&A pairs, covering five modalities. | PMC-15M | 06/2023 |
| ChiMed-VL-Inst [141] | Multiple | 469K | ✓ | QA pairs covering X-rays, CT scans, Echography, etc, translated by GPT-3.5 to Chinese. | PMC-VQA, PMC-CaseReport | 10/2023 |
| CheXInstruct [145] | CXR | 6.1M | ✓ | from existing or curating, five task categories, GPT-4 for translation and formulation. | 65 Datasets | 01/2024 |
| M3D-Data [146] | 3D CT | 662K | ✓ | including tasks of VQA, vision language positioning, and segmentation, using Qwen-72B. | 4 Datasets | 03/2024 |

However, it’s noted that LLM-CXR discards the design of the adapter. It first maps the image to a fixed number of image tokens using pre-trained VQ-GAN [166] based on the reconstructed L2 distance of the two features. Then, these image tokens can be the input or output of LLMs, which is added to the token vocabulary and are random initialized. In this way, LLM-CXR can not only handle the common generation task of CXR-VQA and CXR-to-report generation, but also the report-to-CXR generation, benefiting from VQ-GAN’s ability to generate images.

Beyond handling single modality, some studies are proposed to process multiple image modalities from the perspective of training data construction. For example, BiomedGPT [134] covers pathology, dermatoscope, CT, radiology, and digital camera. Med-PaLM M [138] covers dermatology, mammography, radiograph, pathology, etc. Beyond general 2D images, 3D image process models are also explored for richer spatial information modeling more scalable applications, such as RadFM [139], M3D-LaMed [146], and Dia-LLaMA [147]. Aiming to tackle a wide spectrum of clinical radiology tasks, RadFM is trained on large-scale comprehensive datasets MedMD and RadMD, covering various data modalities (X-ray, CT, MRI, etc), and tasks, featuring over 5000 diseases. It possesses the capability to process both 2D and 3D images. For the 2D images, they are converted into 3D by merely extending an additional dimension. Subsequently, a 3D ViT is employed as the image encoder. Similarly, M3D-LaMed and Dia-LLaMA also employ a 3D ViT as an image encoder. During the tuning process of M3D-LaMed, the 3D image encoder remains frozen, while the proposed 3D spatial pooling perceiver and LLM with LoRA undergo updates. Dia-LLaMA also introduces a disease prototype memory bank as a reference during diagnosis. The disease-aware attention is proposed to extract disease-level representations from visual patches and disease-prototype contrastive loss is used to align these representations with learnable abnormal/normal prototypes. Finally, the predicted diagnostic results can be converted into text prompts using a template description “The {disease name} is [disease state]” to view as an additional input to the LLMs to improve the diagnostic accuracy for infrequent abnormalities.

5.3 Tuning Datasets

To realize medical MLLMs, there are usually two types of data utilized, which can be summarized as alignment and

instruction data. They are concluded in the Table 6.

Alignment Data. It is used for pre-training to align image and text representations. Medical image-text pairs, often found in textbooks or digital libraries, can be converted into alignment data through technical parsing and subsequent processing. For example, PMC-15M [101] is collected by PubMed Parser [167] to process the XML files of PMC and extract captions and the corresponding figure references. After that, items that lack figure references, exhibit syntax errors, or have missing information are systematically eliminated. Upon it, LLaVA-Med-Align [136] is sampled from PMC-15M and ChiMed-VL-Align [141] is similar.

Given the vast number of medical images circulating online, particularly across social media platforms, it would be beneficial and promising to take them into account. OpenPath [22] collects de-identified pathology and their description on Twitter. The 32 pathology-related hashtags are employed to retrieve relevant tweets, and strict protocols are followed regarding inappropriate sample removal and additional text cleaning. Ultimately, 116K image-text pairs from Twitter posts and 59K pairs from the associated replies that received the highest number of likes are retained. CT-RATE [103] comprises chest CT volumes and corresponding radiology text reports from a hospital. It includes about 50K reconstructed CT volumes from 25K distinct CT experiments conducted on 21K unique patients.

Instruction Data. This type of data is used for the instruction tuning. In the medical multimodal domain, it usually composes cross-modal tasks, such as VQA, RG, and coarse/fine-grained image understanding with text. The creation of datasets typically involves using pre-defined prompt templates and answer text, feeding into LLMs to generate questions or to enhance instruction descriptions from a range of original data sources. For example, inspired by the LLaVA-Instruct [130] that process text parts of the multimodal data using LLMs, many studies adopt a similar manner for medical instruction data generation. Biomedical instruction-tuning data LLaVA-Med-Inst [136] is collected using GPT-4, which is filtered from PMC-15M to retain the images that only contain a single plot. Specifically, when presented with an image caption, instructions are formulated in a manner that encourages GPT-4 to produce multi-round questions and answers. This is done in such a way that it appears as if GPT-4 can visualize the image itself, despite the fact that it only has access to the text. CheXInstruct [145] using about 28 public-available CXR datasets for

TABLE 7: Representative benchmark for MLLMs. Columns *LLM* and *H.* indicate using LLM and human experts for evaluation, respectively.

| Benchmark | Modality | Scale | LLM | H. | Tasks (Metrics) | Source | Time |
|---------------------|-----------|-------|-----|----|--|----------------|---------|
| MultiMedBench [138] | Multiple | 1M+ | ✗ | ✗ | QA, RS, VQA, RG, IC (ACC, ROUGE-L, BLEU, F1-RadGraph, CIDEr-D, Macro-AUC, Macro-F1, etc). | 12 Datasets | 07/2023 |
| RedBench [139] | Radiology | 137K | ✗ | ✓ | IC, RG, VQA, rationale diagnosis (ACC, BLEU, ROUGE, UMLS_P, UMLS_R, BERT-Sim, Human, etc). | 13 Datasets | 08/2023 |
| PathQABench [143] | Pathology | 115 | ✗ | ✗ | VQA (ACC, Pathologist evaluation for model comparison). | In-house Cases | 12/2023 |
| CheXbench [145] | CXR | 5K+ | ✓ | ✓ | IC, VQA, RG, RS (ACC, ROUGE-L, CheXbert-S, BERT-S, RadGraph-S, GPT-4, Human). | 7 Datasets | 01/2024 |
| OmniMedVQA [19] | Multiple | 128K | ✗ | ✗ | VQA (ACC). | 73 Datasets | 02/2024 |
| M3D-Bench [146] | 3D CT | 17K+ | ✓ | ✗ | ITR, RG, VQA, positioning, and segmentation (Recall, Qwen-72B, ACC, BERT-Score, IOU, Dice, etc). | 4 Datasets | 03/2024 |

- ◇ Interleaved data: [Text1], [Img1], [Text2], [Img2], [Text3]...
- ◇ Image-text pair: [Img] – [Description]

(a) Interleaved multimodal data and image-text pair.

- ◇ **Coarse-grained Image Understanding:**
Given an [Img], the model is required to diagnose if the [Disease] exists.
- ◇ **Fine-grained Image Understanding:**
Given the [Img], localize the [Region] of [Abnormality].
- ◇ **Text Generation:**
Given the [Img], generate its [Caption].
- ◇ **Question Answering:**
Given the content of the given [Img], answer the [Question].
- ◇ **Miscellaneous:**
Given the [Img], select the text that best matches the image from [Options].

(b) Instruction-following data, from CheXinstruct.

TABLE 8: Multimodal data examples for medical FMs.

generating instruction-tuning datasets. It covers capability, task, dataset, and instance level. It consists of five task categories according to their capabilities: coarse-grained image understanding, fine-grained image understanding, question answering, text generation, and miscellaneous.

In summary, the current trend in datasets for MLLMs involves utilizing advanced LLMs to process text from various sources, such as medical textbooks, digital libraries, or original datasets. The processed data is then transformed into corresponding outputs, ultimately creating image-text alignment data or image-text instruction data. Their illustrations are shown as Table 8.

5.4 Evaluation Benchmarks

To comprehensively explore the capacities of medical MLLMs, several benchmark studies are proposed, such as MultiMedBench [138], RedBench [139], PathQABench [143], CheXbench [145], and M3D-Bench [146]. We give their details in Table 7. They usually integrate multiple datasets of the domain to perform a variety of tasks. For evaluation metrics, three types are commonly utilized, including automatic statistical indicator (e.g., accuracy, ROUGE-L, BLEU, and CIDEr), AI evaluator (e.g., BERT-based and GPT-4), and human expert evaluator.

Typically, MultiMedBench [138] comprises more than 1 million data samples from 12 de-identified datasets of QA, RS (report summarization), VQA, RG, and IC, covering image modalities of pathology, radiograph, genomics, mammography, etc. RadBench [139] encompasses five distinct tasks, including modality recognition, disease diagnosis, VQA, RG, and rationale diagnosis. It has undergone

meticulous manual verification to ensure data quality. It also introduces two additional medical metrics for evaluation, *i.e.*, UMLS_P (precision) and UMLS_R (recall), which aim to measure the overlapping ratio of medical-related words between ground truth and predicted response. The medical-related words are extracted from them by using UMLS. PathQABench [143] and OmniMedVQA [19] utilize the curated VQA datasets for comprehensive evaluation. CheXbench [145] has two evaluation axes, *i.e.*, image perception and textual understanding. The former utilizes six tasks across seven datasets, including view classification, binary disease classification, single disease identification, multi-disease identification, VQA, and image-text reasoning. They are all in the format of multiple-choice and then the accuracy is viewed as the evaluation metric. The latter evaluates the ability of models to generate and summarize text, where a combination of automated metrics (including GPT-4) and human expert evaluations (completeness, correctness, and conciseness) are utilized.

5.5 Medical MLLM Applications

In accordance with general MLLMs, medical MLLMs are typically employed to generate text responses from multimodal input, including medical visual chat, VQA, and RG. For instance, Dia-LLaMA [147] is capable of generating 3D CT reports. LLaVa-Med [136] can handle both medical visual chat and VQA tasks, covering CXR, MRI, histology, gross, and CT domains. The text outputs of MLLMs can also be enhanced by customizing inputs with domain knowledge. To incorporate the expert prior knowledge of radiologists, Kim *et al.* [168] combined the original CXR image with its heatmap that highlights the precise focal points and duration of a radiologist’s attention, which provide extra human intelligence to MLLMs. The experiment results demonstrate performance improvement.

Beyond only generating text responses, medical MLLMs can also give visual image responses, e.g., image synthesis or segmentation with a followed image decoder. Borrowing the capacity of pre-trained VQ-GAN that uses auto-encoding architecture for CXR images, LLM-CXR [133] could generate CXR images using LLM output of predicted virtual images codes by VQ-GAN decoder. MedXChat [144] can implement text-to-CXR synthesis, where the pre-trained stable diffusion (SD) model [169] is used as the foundational framework for CXR generation, which is fine-tuning on MIMIC-CXR dataset using the zero-convolution strategy [170] for the adaptation from general domain to the medical one. The generated prompts by the MLLM are then input into the SD model to generate CXR images. RO-LMM [171] introduces a 3D Residual U-Net [172] for image

processing. A multimodal alignment module is used to integrate comprehensive information from the image encoder and the pre-trained LLM. Subsequently, a 3D image decoder is employed to generate the segmentation mask to realize LMM-assisted breast cancer treatment target segmentation. M3D-LaMed [146] implements referring expression segmentation using a promptable segmentation module, where the last layer embedding of the *[SEG]* token is extracted if it exists in the output. After processing by an MLP, SegVol model [173] is set as the promptable segmentation module to ultimately produce the segmentation mask.

6 DISCUSSIONS OF CURRENT STUDIES

In this section, we will discuss to answer our proposed question from the following five perspectives.

(1). Can the existing multimodal data support advancing intelligent healthcare? Broadly speaking, existing datasets for multimodal healthcare suffer from the issues of diversity, volume, and simplistic construction approaches, which typically hinder the further development of technologies. From §2.3, current RG datasets lack diversity and representation. Medical reports cover a wide variety of conditions, diseases, and clinical scenarios. However, some datasets may be collected in specific fields or healthcare institutions, resulting in samples that are too specific and lack sufficient diversity. Besides, medical reports often rely on rich clinical context. However, existing datasets may not provide sufficient contextual information, such as patient history records and other examination results, which may limit the accuracy and completeness of generation models.

From §2.3, medical VQA datasets encounter the following limitations: 1) Imbalance of quality and quantity. Most datasets are gathered automatically. Despite the significant advancements in the NLP field, achieving 100% accuracy when generating samples remains a challenge. For instance, when three human experts verify samples from the MIMIC-Diff-VQA [45] dataset, the dataset demonstrates an average correctness rate of 97.4% and a minimum correctness rate of 95%. While the VQA-RAD [38] and RadVisDial-Gold [41] datasets are of high quality due to manual collection, their sizes are relatively small, with 3,515 and 500 samples, respectively. Creating a dataset that meets the high standards of both quality and quantity is difficult, given the significant demand for medical professionals. 2) Relative simpleness of the question. Given that the majority of questions are generated automatically based on predefined rules or patterns, the questions in existing medical VQA datasets tend to be simplistic and lacking in variety. For example, *is the lesion associated with a mass effect, what imaging method was used, and what type is the opacity* in VQA-Med-2018 [37], VQA-Med-2019 [39], and MIMIC-Diff-VQA [45], respectively. Although SLAKE [44] incorporates a KG to answer questions that require external medical knowledge, its reasoning schema is in a one-hop direct manner. For wider application prospects, there is a requirement for more comprehensive QA pairs.

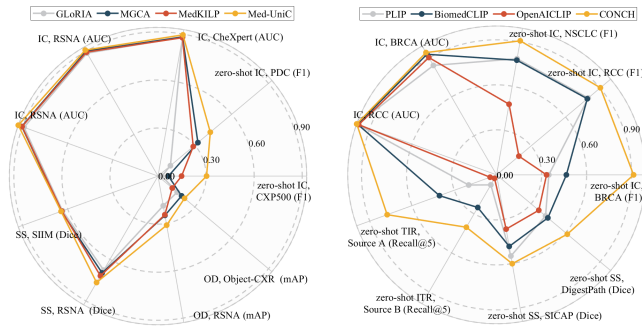
From §5.3, it is observed that the maximum data scale is 16M, which is much less than that in the general domain (e.g., ALIGN [174] with 1.8B image-text pairs and LAION-5B [175]). For instruction data, the samples are generated using prompt engineering with LLMs or collected from

diverse datasets directly, which could lead to simplistic and bias issues. Moreover, there is insufficient emphasis on fine-grained data, which is crucial for medical image perception and implementation of the technologies in reality [62].

(2). Do task-oriented methods effectively address the targeted task? These methods have achieved certain success with the rapid development of AI technologies, but they also confront several significant challenges. Taking RG as an example, existing models have made notable progress in this field. However, these models are still constrained by the data bias inherent in current datasets, which impairs their ability to detect and accurately describe subtle anomalies. On the MIMIC-ABN dataset [176], which exclusively contains anomaly descriptions from the MIMIC-CXR dataset, the performance of these models on various metrics is significantly degraded [73], [176]. For instance, the BLEU-4 and ROUGE scores of RECAP [73] on MIMIC-CXR are 12.5% and 28.8%, respectively, while on MIMIC-ABN, they drop to 8% and 22.3%, respectively. Furthermore, these models often overlook individual variability and the influence of patient-specific clinical context on diagnostic conclusions. This lack of consideration for personalized clinical information limits the overall effectiveness and accuracy of the generated reports. Additionally, while many models perform well on natural language generation metrics such as ROUGE and BLEU, these metrics do not measure clinical accuracy. Consequently, there remains a significant gap between the performance of existing models and the requirements of clinical practice in terms of prediction accuracy, as indicated by metrics like F1-Chexbert. In the field of image generation, existing techniques still suffer from the challenges of quality, accuracy, and interoperability. They often struggle to accurately manipulate specific details, which is critical for medical image analysis [96], [97], [133].

(3). How do FMs contribute to intelligent healthcare? Beyond task-oriented methods that usually undertake one specific task on one modality, FMs contribute to unifying multiple tasks (e.g., IC, ITR, TIR, VQA, and RG) and modalities. Their advantages stem from the utilization of large-scale parameters and training data. Nonetheless, it also introduces challenges, such as complex deployment and diminished efficiency for training and inference [8], [177].

We list some performance results of medical FMs in Figure 6 and Table 9. It can be observed that the CFMs are consistently achieving improvement. But they also need to be fine-tuned for better adapting downstream applications. For instance, IC performance is usually better than the zero-shot IC. Specifically, Med-UniC [111] only achieves about 30% zero-shot IC F1 score on CXP500, and zero-shot TIR, ITR, and SS of CONCH [178] is also relatively low, which would not suffice for practical applications. We also explore the RG ability of MLLMs, it can be seen that MAIRA-1 [142] although achieves competitive results, it even doesn't surpass several non-MLLM models, especially for the clinical metrics that function as a transformation of text that preserves its semantics and prompts the model to concentrate on the report's critical elements without becoming excessively adapted to its style. Beyond that, Hu *et al.* [19] found that medical-specialized MLLMs even exhibit inferior performance to those general-domain models, where BLIP-2 [129] achieves the best performance for all tasks



(a) Results are obtained from Med-UniC [111]. (b) Results are obtained from CONCH [178].

Fig. 6: Performance of medical FMs across various tasks.

TABLE 9: Findings generation performance on the MIMIC-CXR test set. Results are get from MAIRA-1. ‡ means the results are from non-MLLM models.

| Metric | Lexical | | Clinical | | |
|---------|---------|--------|------------------|---------|--------|
| | MAIRA-1 | SOTA | Metric | MAIRA-1 | SOTA |
| ROUGE-L | 28.9 | 27.49 | RadGraph-F1 | 24.3 | 26.71 |
| BLEU-1 | 39.2 | 32.31 | RG _{ER} | 29.6 | 34.7‡ |
| BLEU-4 | 14.2 | 13.30‡ | CheXbert vector | 44.0 | 45.2‡ |
| METEOR | 33.3 | 16.8‡ | RadCliQ (↓) | 3.10 | 3.277‡ |

on average. It indicates that the adaptation of general-purpose MLLMs using existing medical data does not result in emergent properties, underscoring the need for more adaptive and robust MLLMs within the medical domain. In summary, though FMs bring new perspectives to healthcare and possess the ability to undertake multiple tasks, their ability to respond with high precision remains a challenge.

(4). Do the current AI models present any ethical concerns? Despite promising advancements, the integration of AI into healthcare raises significant ethical and regulatory concerns. Issues like data privacy and model bias need careful attention to ensure AI systems are deployed responsibly in this sensitive area. Recent research indicates that FMs might allow data leakage, exposing personal health information when trained on specific data sources [179]. Additionally, biases in FMs often stem from uneven demographic distributions in the training data. As for bias, for example, the study shows that neural models trained on public chest X-ray datasets may underdiagnose in marginalized communities, such as female, black, Hispanic patients, and those insured by medicaid [180]. The broader ethical challenges with FMs, including fairness, accountability, and transparency, are complex and require ongoing attention to minimize ethical risks [1].

Drawing from the superior capabilities of LLMs, medical MLLMs likewise adopt a vulnerability to hallucinations which manifest as irrelevant or factually incorrect responses for the input [181], [182]. This becomes a considerable concern in critical medical situations as there is little room for error. Such errors or hallucinations can raise serious ethical issues for patients or doctors. Different from the general domain, hallucinations in the medical domain can be characterized as multi-tasking, multi-faceted, and hierarchical, making them uniquely challenging and complex to address.

(5). How do professionals assess current multimodal AI technologies? Despite the significant advancements in multimodal AI technologies within healthcare, medical professionals express concerns regarding potential challenges. For example, Hu *et al.* [19] pointed out that despite medical MLLMs claim of robustness, they exhibit inferior performance to those in the general domain, which reveals the limitations inherent in these medical models. Therefore, to develop a more versatile professional, medical MLLMs should continuously incorporate specialized knowledge from various modalities, necessitating significant time and computational resources for refinement. Acosta *et al.* [7] pointed out there existed challenges in curating higher-quality image–text datasets, the exceedingly high number of dimensions contained in multimodal health data, and multimodal model architectures. They highlighted that despite the rapid progress in multimodal learning over recent years, present methods are *unlikely* to be sufficient to address the major challenges.

Based on the above survey and discussions, we can answer the question *has multimodal learning delivered universal intelligence in healthcare?* The answer is **NOT**. The existing research has several significant limitations (from data and technologies to performance and ethics) that need to be addressed to enhance its applicability in practical scenarios.

7 CHALLENGES AND FUTURE DIRECTIONS

Drawing from the advancements in healthcare technology and the aforementioned discussions, we outline the following potential future directions.

(1). High-quality & Diverse Data. In reality, the current intelligent healthcare models are data-centric [183]. However, the datasets employed are somehow simplistic and homogeneous, as described in §6(1). The current success of general LLMs and MLLMs is based on the massive and diverse heterogeneous datasets [1], [18]. To benefit the healthcare community, high-quality and diverse data should be collected for more robust and flexible models applicable to reality. The following aspects could potentially encompass it: effective and varied integration of multiple multimodal datasets, contextual data collection with backgrounds in real-life scenarios [45], [72], user-oriented data construction with multiple image modalities, and domain knowledge alignment data with fine-grained text or images.

(2). Incorporating More Types of Modality. Present models mainly focus on medical images of radiology and pathology. Approaches and techniques for modeling are comparatively well-established and offer the possibility of adaptation to other image modalities [184]. Medical audios (voice recording and stethoscope recording), videos (surgical and patient behavior videos), and time series (ECG, EEG, blood pressure, pulse, and other physiological signal data) can also be incorporated with others or with natural language (*e.g.*, patients’ condition descriptions or domain knowledge [21]) for more comprehensive modeling and more precise medical diagnosis and intervention.

(3). Fine-grained & High-resolution Image Modeling. The medical field necessitates precise visual modeling, given that typically only a minute fraction of the visual features bear relevance to the decision-making process,

while the rest tend to be less informative. Some fine-grained methods should be introduced to concentrate on the local representation of images, rather than only general representation. The potential solution could be the application of multi-granularity, multi-scale, and hierarchical contrastive learning [133], [185]. Further, the model’s capability to process high-resolution images is also crucial. Present FMs mainly focus on low-resolution image processing, *e.g.*, 224×224 [100], [136]. Higher-resolution methods, such as BioViL [106] (512×512) and MAIRA-1 [142] (518×518), demonstrate performance improvements. On one hand, high resolutions would benefit more detailed tissue structure information, clearer delineation of lesion boundaries, and precise quantitative assessment. On the other hand, low-resolution vision encoders cannot directly handle ultra-high resolution medical images, *e.g.*, high-resolution pathology whole-slide images (WSIs) [178]. Some advanced techniques in the general domain can be borrowed, employing high-definition encoders, or independently processing sub-images and subsequently integrating them, which can be referred to Mini-Gemini [186] and Monkey [187].

(4). Effective & Efficient Knowledge Fusion. Vast amounts of domain knowledge in a graph-structured format are stored within meticulously curated knowledge bases, *e.g.*, medical KGs (UMLS and PrimeKG [188]). Evidence confirms that distilling this particular knowledge into designated models is effective [21], [115] and helps to relieve the hallucination of FMs [189]. However, the prevalent methods typically transform structured knowledge into token sequences and feed them into LMs for processing, which would cause information loss. Seamlessly integrating this kind of knowledge into image or language models is challenging due to the inherent representation disparity between their original data. A potential solution lies in learning general semantic tokens [190] for discrete knowledge data, which can be directly incorporated into LMs.

(5). Multimodal In & Multimodal out. Current methods are limited to fixed-form inputs and outputs. Task-oriented methods are usually inputted with one modality and output another modality or prediction. CFMs usually need fine-tuning for downstream tasks and MLLMs concentrate on generating text responses about multimodal inputs. Although Med-Flamingo [137] can be inputted with multiple images, its outputs are fixed to texts. In reality, physicians and patients aspire to establish a diagnosis or predict health conditions based on historical records, where the inputs and outputs would contain multiple medical images or other modalities (*e.g.*, audio and video). For instance, a patient diagnosed with scoliosis may wish to visualize spinal images after a period of treatment. It can be referred to NEX-T-GPT (any-to-any MLLM) [191], utilizing a variety of decoders for any combination outputs of text, image, and video.

(6). Towards a Unified Model. Current models are often designed to process one or a few specific image modalities. However, for broader and more universal applications, it is crucial to develop unified medical models. Given the vast range of medical modalities, organs, diseases, diagnoses, and medications, employing a single-objective model for a specific task proves both challenging and economically inefficient in real-world applications. Integrating primary and various image modalities, as well as both 2D and 3D images,

is an urgent need and has great application prospects. It can be realized from perspectives of data collection [134], [139] or model architecture [184].

(7). Inspiring the Full Potential of FMs. Despite the proven effectiveness of FMs, researchers continue to explore their untapped potential. Early, they discovered that even small differences in similar natural language prompts could significantly impact model performance. This led to various studies on different prompting methods to better utilize FMs, including automatically generated prompts, soft prompts, and pre-trained prompts [192]. Subsequently, the Chain-of-Thought (CoT) [193], [194] method is introduced, where a model explicitly outlines the steps leading to its final answer, improving both transparency and performance. This method gives rise to comparable strategies, like Tree-of-Thought (ToT) [195] and Graph-of-Thought (GoT) [196]. Additionally, retrieval augmentation-based methods have shown promise by incorporating external knowledge to boost FM performance. Overall, these techniques can be tailored for medical FMs, unlocking their full capabilities and enhancing their practicality.

(8). Comprehensive & Unbiased Evaluation Protocol. Beyond the classification tasks that can be directly evaluated by accuracy, current models, especially MLLMs, can generate open-ended answers with free text that is very hard to evaluate. The demographic characteristics of the human evaluators, such as race, sex, age, and other factors, may introduce unconscious biases that affect their assessment of model performance [145]. Automatic statistical metrics like ROUGE-L and BLEU also fail to provide a satisfactory manner. Comprehensive, standardized, and objective evaluation frameworks and metrics are urgent, yet they would be extremely laborious and resource-intensive. Strong AI-based evaluation would provide a promising solution [146].

(9). Enhancing User-oriented Transparency & Interpretability. As an area where there is a strong demand for accountability, the models for medical applications are expected to be able to give explanations of the predictive outcomes that can convince doctors and patients. Although current black-box models can achieve good prediction results, they are not transparent and therefore lack interpretability. For this reason, some explainable reasoning frameworks [197] can be referred to acquire internal logic rules for making predictions. Additionally, with the power of MLLMs, questions can be decomposed into multiple human-understandable symbol subtasks [198] and then solved individually to realize general interpretability.

(10). Minimizing the Risks of Ethics. Some approaches need to be taken to minimize the risks of ethics, aiming to implement robust multimodal architectures [199]. Federated learning offers a solution by training models on local devices, preventing data exposure [200]. Methods like data resampling and model fine-tuning can help align these models with human values to mitigate bias issues [201]. As for the hallucination of MLLMs, it can be mitigated using technologies similar to the general domain, from data, model, training, and inference perspectives [202].

8 CONCLUSION

To explore the open research question *has multimodal learning delivered universal intelligence in healthcare*, we carry out

this study. We first give a comprehensive review of the multimodal datasets, task-oriented techniques, and FMs. Based on them, we carry out the discussion from five sub-questions. Through our study, we discover that the current technologies fall short of realizing the goal of universal intelligent healthcare, highlighting the need for further exploration and development. Additionally, we draw upon our findings to propose ten promising avenues for future research. We anticipate that our study will stimulate extensive discussions among researchers regarding multimodal healthcare, thereby fostering the advancement of this field, especially in the current age of swift advancement of LLMs technologies and their widespread applications.

REFERENCES

- [1] K. He et al. A survey of large language models for healthcare: from data, technology, and applications to accountability and ethics. *arXiv preprint arXiv:2310.05694*, 2023.
- [2] H. Guan and M. Liu. Domain adaptation for medical image analysis: a survey. *IEEE TBME*, 69(3):1173–1185, 2021.
- [3] D. R. Nayak et al. Automated diagnosis of multi-class brain abnormalities using mri images: a deep convolutional neural network based method. *Pattern Recognition Letters*, 138:385–391, 2020.
- [4] A. V. Sadybekov and V. Katritch. Computational approaches streamlining drug discovery. *Nature*, 616(7958):673–685, 2023.
- [5] C.-L. Chi et al. Producing personalized statin treatment plans to optimize clinical outcomes using big data and machine learning. *Journal of biomedical informatics*, 128:104029, 2022.
- [6] K. Singhal et al. Towards expert-level medical question answering with large language models. *arXiv preprint arXiv:2305.09617*, 2023.
- [7] J. N. Acosta et al. Multimodal biomedical ai. *Nature Medicine*, 28(9):1773–1784, 2022.
- [8] P. Shrestha et al. Medical vision language pretraining: A survey. *arXiv preprint arXiv:2312.06224*, 2023.
- [9] Q. Pei et al. Leveraging biomolecule and natural language through multi-modal learning: A survey. *arXiv preprint arXiv:2403.01528*, 2024.
- [10] P. Messina et al. A survey on deep learning and explainability for automatic report generation from medical images. *ACM CSUR*, 54(10s):1–40, 2022.
- [11] X. Wang et al. Tienet: Text-image embedding network for common thorax disease classification and reporting in chest x-rays. In *CVPR*, pp. 9049–9058, 2018.
- [12] T. Baltrusaitis et al. Multimodal machine learning: A survey and taxonomy. *IEEE TPAMI*, 41(2):423–443, 2019.
- [13] K. He et al. Deep residual learning for image recognition. In *CVPR*, pp. 770–778, 2016.
- [14] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, 1997.
- [15] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2016.
- [16] A. Vaswani et al. Attention is all you need. In *NeurIPS*, pp. 5998–6008, 2017.
- [17] A. Radford et al. Learning transferable visual models from natural language supervision. In *ICML*, pp. 8748–8763, 2021.
- [18] S. Yin et al. A survey on multimodal large language models. *CoRR*, abs/2306.13549, 2023.
- [19] Y. Hu et al. Omnimedvqa: A new large-scale comprehensive evaluation benchmark for medical lvlm. In *CVPR*, pp. 22170–22183, 2024.
- [20] Z. Zhao et al. Clip in medical imaging: A comprehensive survey. *arXiv preprint arXiv:2312.07353*, 2023.
- [21] X. Zhang et al. Knowledge-enhanced visual-language pre-training on chest radiology images. *Nature Communications*, 14(1):4542, 2023.
- [22] Z. Huang et al. A visual-language foundation model for pathology image analysis using medical twitter. *Nature Medicine*, 29(9):2307–2316, 2023.
- [23] M. Y. Lu et al. A visual-language foundation model for computational pathology. *Nature Medicine*, pp. 1–12, 2024.
- [24] J. Zhou et al. Pre-trained multimodal large language model enhances dermatological diagnosis using skingpt-4. *Nature Communications*, 15:5649, 2024.
- [25] A. E. Johnson et al. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data*, 6(1):317, 2019.
- [26] O. Bodenreider. The unified medical language system (umls): integrating biomedical terminology. *Nucleic acids research*, 32(suppl_1):D267–D270, 2004.
- [27] D. Demner-Fushman et al. Preparing a collection of radiology examinations for distribution and retrieval. *Journal of the American Medical Informatics Association*, 23(2):304–310, 2016.
- [28] C. Eickhoff et al. Overview of imageclef2017 - image caption prediction and concept detection for biomedical images. In *Working Notes of CLEF 2017*, volume 1866, 2017.
- [29] A. G. S. de Herrera et al. Overview of the imageclef 2018 caption prediction tasks. In *Working Notes of CLEF 2018*, volume 2125, 2018.
- [30] B. Jing et al. On the automatic generation of medical imaging reports. In *ACL*, pp. 2577–2586, 2018.
- [31] O. Pelka et al. Radiology objects in context (ROCO): A multimodal image dataset. In *CVII-STENT and LABELS in MICCAI*, volume 11043, pp. 180–189, 2018.
- [32] A. Bustos et al. Padchest: A large chest x-ray image dataset with multi-label annotated reports. *Medical Image Analysis*, 66:101797, 2020.
- [33] S. Subramanian et al. Medicat: A dataset of medical images, captions, and textual references. In *EMNLP*, pp. 2112–2120, 2020.
- [34] J. Gamper and N. Rajpoot. Multiple instance captioning: Learning representations from histopathology textbooks and articles. In *CVPR*, pp. 16549–16559, 2021.
- [35] M. Li et al. FFA-IR: towards an explainable and reliable medical report generation benchmark. In *NeurIPS*, 2021.
- [36] Y. Tang et al. Work like a doctor: Unifying scan localizer and dynamic generator for automated computed tomography report generation. *Expert Systems with Applications*, 237:121442, 2024.
- [37] S. A. Hasan et al. Overview of imageclef 2018 medical domain visual question answering task. In *Working Notes of CLEF 2018*, volume 2125, 2018.
- [38] J. J. Lau et al. A dataset of clinically generated visual questions and answers about radiology images. *Scientific data*, 5(1):1–10, 2018.
- [39] A. Ben Abacha et al. Vqa-med: Overview of the medical visual question answering task at imageclef 2019. In *Proceedings of CLEF*, 2019.
- [40] A. B. Abacha et al. Overview of the vqa-med task at imageclef 2020: Visual question answering and generation in the medical domain. In *Working Notes of CLEF 2020*, volume 2696, 2020.
- [41] O. Kovaleva et al. Towards visual dialog for radiology. In *Proceedings of the 19th SIGBioMed Workshop*, pp. 60–69, 2020.
- [42] X. He et al. Pathvqa: 30000+ questions for medical visual question answering. *arXiv preprint arXiv:2003.10286*, 2020.
- [43] A. Ben Abacha et al. Overview of the vqa-med task at imageclef 2021: Visual question answering and generation in the medical domain. In *Proceedings of the CLEF 2021*, 2021.
- [44] B. Liu et al. SLAKE: A semantically-labeled knowledge-enhanced dataset for medical visual question answering. In *ISBI*, pp. 1650–1654. IEEE, 2021.
- [45] X. Hu et al. Expert knowledge-aware image difference graph representation learning for difference-aware medical visual question answering. In *KDD*, pp. 4156–4165, 2023.
- [46] Y.-D. Zhang et al. Advances in multimodal data fusion in neuroimaging: Overview, challenges, and novel orientation. *Information Fusion*, 64:149–187, 2020.
- [47] S. Li et al. Pixel-level image fusion: A survey of the state of the art. *information Fusion*, 33:100–112, 2017.
- [48] D. S. Shibu and S. S. Priyadharsini. Multi scale decomposition based medical image fusion using convolutional neural network and sparse representation. *Biomedical Signal Processing and Control*, 69:102789, 2021.
- [49] Y. Liu et al. A general framework for image fusion based on multi-scale transform and sparse representation. *Information fusion*, 24:147–164, 2015.
- [50] V. Bhavana and H. Krishnappa. Multi-modality medical image fusion using discrete wavelet transform. *Procedia Computer Science*, 70:625–631, 2015.

- [51] L. Huang et al. Deep evidential fusion with uncertainty quantification and contextual discounting for multimodal medical image segmentation. *arXiv preprint arXiv:2309.05919*, 2023.
- [52] M. Safari et al. Medfusiongan: multimodal medical image fusion using an unsupervised deep generative adversarial network. *BMC Medical Imaging*, 23(1):203, 2023.
- [53] S. Zhang et al. Multi-modal contrastive mutual learning and pseudo-label re-learning for semi-supervised medical image segmentation. *Medical Image Analysis*, 83:102656, 2023.
- [54] Y. Liu et al. Sf-net: A multi-task model for brain tumor segmentation in multimodal mri via image fusion. *IEEE Signal Processing Letters*, 29:1799–1803, 2022.
- [55] A. Dosovitskiy et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- [56] G. Piella. A general framework for multiresolution image fusion: from pixels to regions. *Information fusion*, 4(4):259–280, 2003.
- [57] S. U. R. Khan et al. Hybrid-net: A fusion of densenet169 and advanced machine learning classifiers for enhanced brain tumor diagnosis. *International Journal of Imaging Systems and Technology*, 34(1):e22975, 2024.
- [58] P. H. Foo and G. W. Ng. High-level information fusion: An overview. *J. Adv. Inf. Fusion*, 8(1):33–72, 2013.
- [59] I. Najdenkoska et al. Variational topic inference for chest x-ray report generation. In *MICCAI*, volume 12903, pp. 625–635, 2021.
- [60] Y. Li et al. Unify, align and refine: Multi-level semantic alignment for radiology report generation. In *ICCV*, pp. 2863–2874, 2023.
- [61] Y. Chen et al. Representative image feature extraction via contrastive learning pretraining for chest x-ray report generation. *CoRR*, abs/2209.01604, 2022.
- [62] S. Wang et al. Fine-grained medical vision-language representation learning for radiology report generation. In *EMNLP*, pp. 15949–15956, 2023.
- [63] B. Jing et al. Show, describe and conclude: On exploiting the structure information of chest x-ray reports. In *ACL*, pp. 6570–6580, 2019.
- [64] J. Delbrouck et al. Improving the factual correctness of radiology report generation with semantic rewards. In *Findings of EMNLP*, pp. 4348–4360, 2022.
- [65] S. Jain et al. Radgraph: Extracting clinical entities and relations from radiology reports. In *NeurIPS*, 2021.
- [66] D. Parres et al. Improving radiology report generation quality and diversity through reinforcement learning and text augmentation. *Bioengineering*, 11(4):351, 2024.
- [67] T. Zhang et al. Bertscore: Evaluating text generation with BERT. In *ICLR*, 2020.
- [68] B. Jing et al. On the automatic generation of medical imaging reports. In *ACL*, pp. 2577–2586, 2018.
- [69] M. Li et al. Dynamic graph enhanced contrastive learning for chest x-ray report generation. In *CVPR*, pp. 3334–3343, 2023.
- [70] Z. Huang et al. Kiut: Knowledge-injected u-transformer for radiology report generation. In *CVPR*, pp. 19809–19818, 2023.
- [71] W. Hou et al. Organ: Observation-guided radiology report generation via tree reasoning. In *ACL*, pp. 8108–8122, 2023.
- [72] S. Bannur et al. Learning to exploit temporal structure for biomedical vision-language processing. In *CVPR*, pp. 15016–15027, 2023.
- [73] W. Hou et al. RECAP: towards precise radiology report generation via dynamic disease progression reasoning. In *Findings of EMNLP*, pp. 2134–2147, 2023.
- [74] Y. Huang et al. OVQA: A clinically generated visual question answering dataset. In *SIGIR*, pp. 2924–2938, 2022.
- [75] B. D. Nguyen et al. Overcoming data limitation in medical visual question answering. In *MICCAI*, volume 11767, pp. 522–530, 2019.
- [76] T. Do et al. Multiple meta-model quantifying for medical visual question answering. In *MICCAI*, volume 12905, pp. 64–74, 2021.
- [77] M. Wang et al. Medical visual question answering based on question-type reasoning and semantic space constraint. *Artificial Intelligence in Medicine*, 131:102346, 2022.
- [78] J. Devlin et al. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, pp. 4171–4186, 2019.
- [79] Z. Chen et al. Align, reason and learn: Enhancing medical vision-and-language pre-training with knowledge. In *ACM MM*, pp. 5152–5161, 2022.
- [80] J. Liu et al. Parameter-efficient transfer learning for medical visual question answering. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2023.
- [81] X. Liu et al. Generating binary tags for fast medical image retrieval based on convolutional nets and radon transform. In *IJCNN*, pp. 2872–2878, 2016.
- [82] A. Mbilyi and H. Schuldt. Cross-modality medical image retrieval with deep features. In *BIBM*, pp. 2632–2639, 2020.
- [83] L. Xu et al. Multi-manifold deep discriminative cross-modal hashing for medical image retrieval. *IEEE TIP*, 31:3371–3385, 2022.
- [84] Y. Zhang et al. Category supervised cross-modal hashing retrieval for chest x-ray and radiology reports. *Computers & Electrical Engineering*, 98:107673, 2022.
- [85] G. Ding et al. Semantic extension for cross-modal retrieval of medical image-diagnosis report. In *NLPCC*, pp. 442–455, 2023.
- [86] Z. Li et al. Lvit: Language meets vision transformer in medical image segmentation. *IEEE TMI*, 2023.
- [87] Y. Zhao et al. Dtan: Diffusion-based text attention network for medical image segmentation. *Computers in Biology and Medicine*, 168:107728, 2024.
- [88] Z. Dong et al. Diffusion model-based text-guided enhancement network for medical image segmentation. *Expert Systems with Applications*, pp. 123549, 2024.
- [89] I. Goodfellow et al. Generative adversarial nets. volume 27, 2014.
- [90] D. Nie et al. Medical image synthesis with context-aware generative adversarial networks. In *MICCAI*, pp. 417–425, 2017.
- [91] Y. Hiasa et al. Cross-modality image synthesis from unpaired data using cyclegan: Effects of gradient consistency loss and training data size. In *SASHIMI*, pp. 31–41, 2018.
- [92] A. Ben-Cohen et al. Cross-modality synthesis from ct to pet using fcn and gan networks for improved automated lesion detection. *Engineering Applications of Artificial Intelligence*, 78:186–194, 2019.
- [93] Y. Pan et al. Synthesizing missing pet from mri with cycle-consistent generative adversarial networks for alzheimer’s disease diagnosis. In *MICCAI*, pp. 455–463, 2018.
- [94] H. Choi and D. S. Lee. Generation of structural mr images from amyloid pet: application to mr-less quantification. *Journal of Nuclear Medicine*, 59(7):1111–1117, 2018.
- [95] J. Ho et al. Denoising diffusion probabilistic models. In *NeurIPS*, volume 33, pp. 6840–6851, 2020.
- [96] Q. Lyu and G. Wang. Conversion between ct and mri images using diffusion and score-matching models. *arXiv preprint arXiv:2209.12104*, 2022.
- [97] X. Meng et al. A novel unified conditional score-based generative framework for multi-modal medical image completion. *arXiv preprint arXiv:2207.03430*, 2022.
- [98] Y. Zhang et al. Contrastive learning of medical visual representations from paired images and text. In *MLHC*, volume 182, pp. 2–25, 2022.
- [99] S. Esлами et al. Pubmedclip: How much does clip benefit visual question answering in the medical domain? In *Findings of EACL*, pp. 1181–1193, 2023.
- [100] E. Tiu et al. Expert-level detection of pathologies from unannotated chest x-ray images via self-supervised learning. *Nature Biomedical Engineering*, 6(12):1399–1406, 2022.
- [101] S. Zhang et al. Biomedclip: a multimodal biomedical foundation model pretrained from fifteen million scientific image-text pairs. *arXiv preprint arXiv:2303.00915*, 2023.
- [102] Y. Sun et al. Pathasst: A generative foundation ai assistant towards artificial general intelligence of pathology. In *AAAI*, volume 38, pp. 5034–5042, 2024.
- [103] I. E. Hamamci et al. A foundation model utilizing chest ct volumes and radiology reports for supervised-level zero-shot detection of abnormalities. *arXiv preprint arXiv:2403.17834*, 2024.
- [104] Y. Xie et al. Pairaug: What can augmented image-text pairs do for radiology? In *CVPR*, pp. 11652–11661, 2024.
- [105] S. Huang et al. Gloria: A multimodal global-local representation learning framework for label-efficient medical image recognition. In *ICCV*, pp. 3922–3931, 2021.
- [106] B. Boecking et al. Making the most of text semantics to improve biomedical vision-language processing. In *ECCV*, pp. 1–21, 2022.
- [107] Z. Wang et al. Medclip: Contrastive learning from unpaired medical images and text. In *EMNLP*, pp. 3876–3887, 2022.
- [108] F. Wang et al. Multi-granularity cross-modal alignment for generalized medical visual representation learning. In *NeurIPS*, 2022.
- [109] C. Wu et al. Medklip: Medical knowledge enhanced language-image pre-training for x-ray diagnosis. In *ICCV*, pp. 21315–21326. IEEE, 2023.

- [110] Z. Chen et al. Towards unifying medical vision-and-language pre-training via soft prompts. In *ICCV*, pp. 23403–23413, 2023.
- [111] Z. Wan et al. Med-unic: Unifying cross-lingual medical vision-language pre-training by diminishing bias. *NeurIPS*, 36, 2023.
- [112] Z. Wei et al. Masked contrastive reconstruction for cross-modal medical image-report retrieval. *arXiv preprint arXiv:2312.15840*, 2023.
- [113] Z. Li et al. Mlip: Enhancing medical visual representation with divergence encoder and knowledge-guided contrastive learning. *arXiv preprint arXiv:2402.02045*, 2024.
- [114] M. H. Phan et al. Decomposing disease descriptions for enhanced pathology detection: A multi-aspect vision-language matching framework. *arXiv preprint arXiv:2403.07636*, 2024.
- [115] X. Zhou et al. Knowledge-enhanced visual-language pretraining for computational pathology. *arXiv preprint arXiv:2404.09942*, 2024.
- [116] H. Luo et al. Devide: Faceted medical knowledge for improved medical vision-language pre-training. *arXiv preprint arXiv:2404.03618*, 2024.
- [117] M. Moor et al. Foundation models for generalist medical artificial intelligence. *Nature*, 616(7956):259–265, 2023.
- [118] K. He et al. Deep residual learning for image recognition. In *CVPR*, pp. 770–778, 2016.
- [119] Y. Liu et al. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [120] E. Alsentzer et al. Publicly available clinical bert embeddings. *arXiv preprint arXiv:1904.03323*, 2019.
- [121] Y. Gu et al. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare*, 3(1):1–23, 2021.
- [122] X. Zhou et al. Adaptive sparse pairwise loss for object re-identification. In *CVPR*, pp. 19691–19701, 2023.
- [123] A. Q. Jiang et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- [124] O. Ronneberger et al. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, pp. 234–241, 2015.
- [125] F. I. Diakogiannis et al. Resunet-a: A deep learning framework for semantic segmentation of remotely sensed data. *ISPRS Journal of Photogrammetry and Remote Sensing*, 162:94–114, 2020.
- [126] J. Redmon and A. Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.
- [127] J.-B. Alayrac et al. Flamingo: a visual language model for few-shot learning. *NeurIPS*, 35:23716–23736, 2022.
- [128] J. Yu et al. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*, 2022.
- [129] J. Li et al. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, pp. 19730–19742, 2023.
- [130] H. Liu et al. Visual instruction tuning. In *NeurIPS*, 2023.
- [131] D. Zhu et al. Minigpt-4: Enhancing vision-language understanding with advanced large language models. In *ICLR*, 2024.
- [132] Q. Chen et al. Medblip: Bootstrapping language-image pre-training from 3d medical images and texts. *arXiv preprint arXiv:2305.10799*, 2023.
- [133] S. Lee et al. Llm-cxr: Instruction-finetuned llm for cxr image understanding and generation. In *ICLR*, 2024.
- [134] K. Zhang et al. Biomedgpt: A unified and generalist biomedical generative pre-trained transformer for vision, language, and multimodal tasks. *arXiv preprint arXiv:2305.17100*, 2023.
- [135] O. Thawkar et al. Xraygpt: Chest radiographs summarization using medical vision-language models. *arXiv preprint arXiv:2306.07971*, 2023.
- [136] C. Li et al. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. *NeurIPS*, 36, 2023.
- [137] M. Moor et al. Med-flamingo: a multimodal medical few-shot learner. In *Machine Learning for Health (ML4H)*, pp. 353–367, 2023.
- [138] T. Tu et al. Towards generalist biomedical ai. *NEJM AI*, 1(3):A10a2300138, 2024.
- [139] C. Wu et al. Towards generalist foundation model for radiology. *arXiv preprint arXiv:2308.02463*, 2023.
- [140] C. Pellegrini et al. Radialog: A large vision-language model for radiology report generation and conversational assistance. *arXiv preprint arXiv:2311.18681*, 2023.
- [141] J. Liu et al. Qilin-med-vl: Towards chinese large vision-language model for general healthcare. *arXiv preprint arXiv:2310.17956*, 2023.
- [142] S. L. Hyland et al. Maira-1: A specialised large multimodal model for radiology report generation. *arXiv preprint arXiv:2311.13668*, 2023.
- [143] M. Y. Lu et al. A foundational multimodal vision language ai assistant for human pathology. *arXiv preprint arXiv:2312.07814*, 2023.
- [144] L. Yang et al. Medxchat: Bridging cxr modalities with a unified multimodal large model. *arXiv preprint arXiv:2312.02233*, 2023.
- [145] Z. Chen et al. Chexagent: Towards a foundation model for chest x-ray interpretation. *arXiv preprint arXiv:2401.12208*, 2024.
- [146] F. Bai et al. M3d: Advancing 3d medical image analysis with multi-modal large language models. *arXiv preprint arXiv:2404.00578*, 2024.
- [147] Z. Chen et al. Dia-llama: Towards large language model-driven ct report generation. *arXiv preprint arXiv:2403.16386*, 2024.
- [148] J. M. Zambrano Chaves et al. Training small multimodal models to bridge biomedical competency gap: A case study in radiology imaging. *arXiv e-prints*, pp. arXiv–2403, 2024.
- [149] S. Kang et al. Wolf: Large language model framework for cxr understanding. *arXiv preprint arXiv:2403.15456*, 2024.
- [150] A. Brock et al. High-performance large-scale image recognition without normalization. In *ICML*, pp. 1059–1071. PMLR, 2021.
- [151] Y. Fang et al. Eva: Exploring the limits of masked visual representation learning at scale. In *CVPR*, pp. 19358–19369, 2023.
- [152] F. Pérez-García et al. Rad-dino: Exploring scalable medical image encoders beyond text supervision. *arXiv preprint arXiv:2401.10815*, 2024.
- [153] M. Caron et al. Emerging properties in self-supervised vision transformers. In *ICCV*, pp. 9650–9660, 2021.
- [154] S. Zhang et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- [155] H. W. Chung et al. Scaling instruction-finetuned language models. *arXiv e-prints*, pp. arXiv–2210, 2022.
- [156] W.-L. Chiang et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023.
- [157] A. Q. Jiang et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- [158] H. Touvron et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [159] H. Touvron et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [160] Y. Cui et al. Efficient and effective text encoding for chinese llama and alpaca. *arXiv preprint arXiv:2304.08177*, 2023.
- [161] R. Luo et al. Biogpt: generative pre-trained transformer for biomedical text generation and mining. *Briefings in bioinformatics*, 23(6):bbac409, 2022.
- [162] A. Venigalla et al. Biomedlm: a domain-specific large language model for biomedical text. *MosaicML*. Accessed: Dec, 23(3):2, 2022.
- [163] C. Wu et al. Pmc-llama: Further finetuning llama on medical papers. *arXiv preprint arXiv:2304.14454*, 2023.
- [164] E. J. Hu et al. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [165] P. Lu et al. Learn to explain: Multimodal reasoning via thought chains for science question answering. *NeurIPS*, 35:2507–2521, 2022.
- [166] P. Esser et al. Taming transformers for high-resolution image synthesis. In *CVPR*, pp. 12873–12883, 2021.
- [167] T. Achakulvisut et al. Pubmed parser: A python parser for pubmed open-access xml subset and medline xml dataset xml dataset. *Journal of Open Source Software*, 5(46):1979, 2020.
- [168] Y. Kim et al. Enhancing human-computer interaction in chest x-ray analysis using vision and language model with eye gaze patterns. *arXiv preprint arXiv:2404.02370*, 2024.
- [169] R. Rombach et al. High-resolution image synthesis with latent diffusion models. In *CVPR*, pp. 10684–10695, 2022.
- [170] L. Zhang et al. Adding conditional control to text-to-image diffusion models. In *ICCV*, pp. 3836–3847, 2023.
- [171] K. Kim et al. Lmm-assisted breast cancer treatment target segmentation with consistency embedding. *arXiv preprint arXiv:2311.15876v2*, 2023.
- [172] Ö. Çiçek et al. 3d u-net: learning dense volumetric segmentation from sparse annotation. In *MICCAI*, pp. 424–432, 2016.
- [173] Y. Du et al. Segvol: Universal and interactive volumetric medical image segmentation. *arXiv preprint arXiv:2311.13385*, 2023.
- [174] C. Jia et al. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, volume 139, pp. 4904–4916, 2021.

- [175] C. Schuhmann et al. LAION-5B: an open large-scale dataset for training next generation image-text models. In *NeurIPS*, 2022.
- [176] J. Ni et al. Learning visual-semantic embeddings for reporting abnormal findings on chest x-rays. In *Findings of EMNLP*, pp. 1954–1960, 2020.
- [177] B. Azad et al. Foundational models in medical imaging: A comprehensive survey and future vision. *arXiv preprint arXiv:2310.18689*, 2023.
- [178] R. J. Chen et al. Towards a general-purpose foundation model for computational pathology. *Nature Medicine*, 30(3):850–862, 2024.
- [179] J. Huang et al. Are large pre-trained language models leaking your personal information?, 2022.
- [180] L. Seyyed-Kalantari et al. Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations. *Nature medicine*, 27(12):2176–2182, 2021.
- [181] J. Chen et al. Detecting and evaluating medical hallucinations in large vision language models. *arXiv preprint arXiv:2406.10185*, 2024.
- [182] A. Pal and M. Sankarasubbu. Gemini goes to med school: exploring the capabilities of multimodal large language models on medical challenge problems & hallucinations. *arXiv preprint arXiv:2402.07023*, 2024.
- [183] Y. Zhang et al. Data-centric foundation models in computational healthcare: A survey. *arXiv preprint arXiv:2401.02458*, 2024.
- [184] Y. Ye et al. Continual self-supervised learning: Towards universal multi-modal medical data representation learning. 2024.
- [185] Q. Lin et al. Contrastive graph representations for logical formulas embedding. *IEEE TKDE*, 35(4):3563–3574, 2023.
- [186] Y. Li et al. Mini-gemini: Mining the potential of multi-modality vision language models. *arXiv preprint arXiv:2403.18814*, 2024.
- [187] Z. Li et al. Monkey: Image resolution and text label are important things for large multi-modal models. In *CVPR*, pp. 26763–26773, 2024.
- [188] P. Chandak et al. Building a knowledge graph to enable precision medicine. *Scientific Data*, 10(1):67, 2023.
- [189] X. Guan et al. Mitigating large language model hallucinations via autonomous knowledge graph-based retrofitting. In *AAAI*, volume 38, pp. 18126–18134, 2024.
- [190] Z. Liu et al. Rethinking tokenizer and decoder in masked graph modeling for molecules. *NeurIPS*, 36, 2023.
- [191] S. Wu et al. Next-gpt: Any-to-any multimodal llm. In *ICML*, 2024.
- [192] T. Shin et al. Autoprompt: Eliciting knowledge from language models with automatically generated prompts. In *EMNLP*, pp. 4222–4235, 2020.
- [193] J. Wei et al. Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*, volume 35, pp. 24824–24837, 2022.
- [194] Z. Zhang et al. Multimodal chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*, 2023.
- [195] S. Yao et al. Tree of thoughts: Deliberate problem solving with large language models. In *NeurIPS*, volume 36, 2023.
- [196] M. Besta et al. Graph of thoughts: Solving elaborate problems with large language models. In *AAAI*, volume 38, pp. 17682–17690, 2024.
- [197] Q. Lin et al. Techs: Temporal logical graph networks for explainable extrapolation reasoning. In *ACL*, pp. 1281–1293, 2023.
- [198] F. Xu et al. Symbol-llm: Towards foundational symbol-centric interface for large language models. *ACL*, 2024.
- [199] J. Ma et al. Robust visual question answering: Datasets, methods, and future challenges. *IEEE TPAMI*, 2024.
- [200] L. Li et al. A review of applications in federated learning. *Computers & Industrial Engineering*, 149:106854, 2020.
- [201] J. Schulman et al. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [202] Z. Bai et al. Hallucination of multimodal large language models: A survey. *arXiv preprint arXiv:2404.18930*, 2024.

APPENDIX

A. DATASET DETAILS

A.1 Datasets for Report Generation

There are several common datasets for report generation, as it is a very significant task where experts can communicate using language texts. We summarize these studies in Table 2. We can observe that they mainly focus on the radiology

modality. Current report generation datasets are mainly constructed using digital databases and libraries, e.g., PubMed.

At the beginning, IU X-ray [27] undergoes an automated de-identification process for X-ray images and their responding reports, which is subsequently validated manually to ensure accuracy. Each report consists of five sections: *indications, findings, impression, manual encoding*, and *MTI encoding*. The former three denote the symptoms or reasons for examination, radiology observations, and final diagnosis, respectively. The results of the manual and automatic coding of findings are embedded in the latter two. ICLEF-Caption-2017 [28] collect 184.6K image-caption pairs from scholarly biomedical articles in PMC, which is filtered from 3 million images automatically. However, due to the fully automated nature of the construction method, the dataset contains a certain level of noise, notably 10~20% of the images are either compound or non-clinical. Based on ICLEF-Caption-2017, ICLEF-Caption-2018 [29] further improves the quality and quantity of the dataset, which utilizes CNNs to decline noise and finally there are 232.3K image-caption pairs. To explore the interaction between visual elements and semantic relationships inherent in radiology images, Pelka *et al.* [31] automatically construct ROCO which is a fine-grained multimodal image dataset. This intricately detailed multimodal dataset encompasses a variety of medical imaging modalities, such as Computer Tomography, Ultrasound, X-ray, Fluoroscopy, Positron Emission Tomography, Mammography, Magnetic Resonance Imaging, and Angiography. It is filtered using pre-trained neural networks, and the final text is linked to UMLS CUIs and semantic types for data standardization and additional image interrelations. PadChest [32] is a large-scale and high-resolution CXR dataset. For data gathering, 27% of the reports are meticulously annotated by skilled physicians, while the rest is automatically labeled utilizing a supervised methodology that leverages the capabilities of an RNN with attention mechanisms. Specifically, the reports are categorized using 174 distinct radiographic findings, 19 differential diagnoses, and 104 anatomical locations. These categories are structured as a hierarchical taxonomy and correspond to the standardized terminology of the UMLS.

For the pathology images, PEIR Gross [30] collects images and their description captions from the Pathology Education Informational Resource (PEIR) digital library. It only utilizes gross lesion images from 21 PEIR pathology sub-categories and each caption contains only one sentence that has an average of 12.0 words. To enhance the dense supervision of tasks related to computational pathology, the ARCH dataset [34] has been introduced. It possesses a multiple instance feature, where a specific caption could be related to multiple images, referred to as a *bag*. The image-caption pairs are meticulously curated from PubMed articles and 10 different textbooks, followed by a comprehensive filtration process.

Given that medical figures in particular are typically complex and often comprised of multiple subfigures, MediCaT [33] incorporates subcaptions and inline references to facilitate a comprehensive understanding of the relationship between text and figures, which constructs a dataset of medical images in context. Image inline references are typically located in the main body of the paper, significantly

distanced from the pertinent figure, which often leads to them being overlooked. By employing MedICaT, research can be conducted to correlate each subfigure with its corresponding subcaption within a compound figure. MedICaT surpasses previous datasets in its comprehensiveness, as it includes inline references for 74% of the figures within the dataset. Different from other main datasets that need to filter out the compound figures, ARCH and MedICaT are allowed compound figures.

To enhance the trustworthiness of the diagnostic methods as existing methods can only predict reports without accurate explanation, FFA-IR [35] includes annotations of 46 categories of lesions with a total of 12,166 regions along with FFA images. The annotation schema was developed by the ophthalmologists based on their expert knowledge, and covers most typical retinal lesions. 3D CT report dataset CTRG [36], including brain and chest splits. Subjective descriptions unrelated to the diagnosis are filtered out of the report. The report is then split into a template and abnormal descriptions to facilitate model training.

A.2 Datasets for Medical VQA

As medical images are typically accompanied by corresponding captions in the database, creating medical VQA datasets primarily involves selecting images, choosing sentences from the captions as answers, and finally generating questions. We summarize the main characteristics of the current common-used dataset for medical VQA in Table 3.

As a groundbreaking study, VQA-Med-2018 [37] introduces a semi-automatic approach to generate question-answer pairs from captions of the radiology images sourced from PubMed, released on the ImageCLEF platform⁴ (based on ICLEF-Caption-2017). A rule-based system is employed to generate possible question-answer pairs from captions and then the question ranking process is conducted to select the most appropriate question. Finally, two expert human annotators carefully review all the generated samples to filter out noisy and incorrect ones. The text similarity metrics BLEU as well as word-based and concept-based semantic similarity are introduced for evaluation. Similarly to VQA-Med-2018, VQA-Med-2019 [39], VQA-Med-2020 [40], and VQA-Med-2021 [43] automatically generate radiology VQA samples based on MedPix database. They all generate specific questions about sentences in image captions using predefined rules or patterns. It is noted that VQA-Med-2019 also utilize the accuracy metric of exact matching for generation evaluation, emphasizing the precise match between a predicted answer and the ground truth answer. To guarantee the quality and investigate the question in a realistic scene, VQA-RAD [38] dataset is carried out, which is the first manually constructed dataset. Clinicians are naturally presented with questions regarding radiology images (CT, MRI, X-ray) and are expected to provide corresponding answers.

Different from the above generation datasets, the answer format in RadVisDial-Silver and RadVisDial-Gold [41] is structured as a multiple-choice selection among four options. It indicates that they involve classification tasks and the question contents are from 13

abnormalities. RadVisDial-Silver and RadVisDial-Gold are gathered through automated and manual means, respectively. PathVQA [42] is the first attempt to build a VQA dataset about pathology images. It extracts pathology images and their captions from both electronic pathology textbooks and the Pathology Education Informational Resource (PEIR) Digital Library⁵. Subsequently, sentence simplification, question transducer, and post-processing are introduced to generate questions. SLAKE [44], a bilingual dataset that contains samples in both English and Chinese, emphasizes the semantic labels and a structured medical knowledge graph in its features. It offers two types of semantic visual annotations for each radiology image: masks for semantic segmentation and bounding boxes for object detection. In addition to basic clinical questions, it introduces compositional questions that demand multiple reasoning steps utilizing a personalized medical KG derived from OwnThink⁶. To illustrate the difference between the two medical images for monitoring changes in patients' physical conditions, the MIMIC-Diff-VQA [45] dataset is introduced. The initial step involves gathering keywords, such as names of abnormalities and important attributes like location, level, and type, from the MIMIC-CXR dataset. Next, an intermediate KeyInfo database is created, and study pair questions are generated based on specific patterns. Each pair of images consists of a main image and a reference image, derived from distinct studies of the identical patient. The selection of main and reference visits adheres strictly to the initial visit being designated as the *reference* and the subsequent visit as the *main* image. To facilitate an all-inclusive assessment of large vision-language models within the medical domain, Hu *et al.* [19] construct OmniMedVQA dataset. The initial phase of this process involves the collection of 75 distinct medical classification datasets encompassing 12 fine-grained imaging modalities, covering over 20 human anatomical regions. Subsequently, unique QA templates are devised for each dataset to generate questions based on the original category labels. The final step refines the QA pairs using GPT-3.5 and human validation checks.

4. <https://www.imageclef.org/>

5. <http://peir.path.uab.edu/library/index.php?/category/2>

6. <https://www.ownthink.com>