# Detecting Any Instruction-to-Answer Interaction Relationship:Universal Instruction-to-Answer Navigator for Med-VQA

Zhongze Wu [1] [*]  Hongyan Xu [2] [*]  Yitian Long [3]  Shan You [4]  Xiu Su [1 5] [†]  Jun Long [1] [†]  Yueyi Luo [1]  Chang Xu [5]

## Abstract

Medical Visual Question Answering (Med-VQA) interprets complex medical imagery using user instructions for precise diagnostics, yet faces challenges due to diverse, inadequately annotated images. In this paper, we introduce the Universal Instruction-to-Answer Navigator (Uni-Med) framework for extracting instruction-to-answer relationships, facilitating the understanding of visual evidence behind responses. Specifically, we design the Instruct-to-Answer Clues Interpreter (IAI) to generate visual explanations based on the answers and mark the core part of instructions with "real intent" labels. The IAI-Med VQA dataset, produced using IAI, is now publicly available to advance Med-VQA research. Additionally, our Token-Level Cut-Mix module dynamically aligns visual explanations with image patches, ensuring answers are traceable and learnable. We also implement intention-guided attention to minimize non-core instruction interference, sharpening focus on 'real intent'. Extensive experiments on SLAKE datasets show Uni-Med's superior accuracies (87.52% closed, 86.12% overall), outperforming MedVInT-PMC-VQA by 1.22% and 0.92%.

## 1. Introduction

With the ongoing development of artificial intelligence (Su et al., 2022b; 2021a; Li et al., 2023c), we have observed significant potential of large model technologies, in downstream applications (Xu et al., 2022; Su et al., 2021b).

[*]Equal contribution [1]Central South University, Changsha, Hunan, China [2]University of New South Wales, Sydney, Australia [3]Vanderbilt University, Nashville, Tennessee, USA [4]SenseTime [5]University of Sydney, Sydney, Australia. Correspondence to: Xiu Su <xisu5992@uni.sydney.edu.au>, Jun Long <junlong@csu.edu.cn>.

Specifically, the advancements in multi-modal large language models (MLLMs), evidenced by studies like (Li et al., 2023a; Moor et al., 2023), have demonstrated significant potential in biomedical applications. Among these, the introduction of models like GPT-4V (OpenAI, 2023) has contributed to these improvements, particularly in medical diagnostics, as seen in works by (Li et al., 2023d), (Han et al., 2023) and (Wang et al., 2023). These models aid in processing comprehensive multi-modal information (Cao et al., 2023) and reducing language ambiguities , beneficial especially for non-experts.

Despite these advancements, effective application of MLLMs(Liu et al., 2023) in biomedical vision tasks, especially in data-scarce environments, remains challenging. Issues like an over-reliance (Su et al., 2021c) on text and image labels in medical diagnosis and analysis errors continue to be prevalent in current models (Yan et al., 2023). Recent works focus on designing various prompts (Chen et al., 2023; Zhan et al., 2023; Cao et al., 2024) to reduce modality interference and enhance following ability to instructions. However, these models often face difficulties with modal interference, impacting their capacity to accurately comprehend user instructions and deduce answers (Ye et al., 2023).

Current pretraining approaches, such as those in (Wang et al., 2022; Li et al., 2023a; Zhang et al., 2023b; Li et al., 2023b), face limitations in providing fine-grained explanations and effectively following instructions for medical queries (Chen et al., 2024; Cong et al., 2022). The post-hoc attention method by (Zhang et al., 2023a) enhances instruction-following capabilities but lacks automated, user-intent-focused training. Moreover, while various explicit data enhancement techniques (Tang et al., 2020a; Guo, 2020; Gong et al., 2022a) have been developed, they lack in providing dynamic, instruction-specific feature-level enhancements during training, leading to potential misalignments between learned representations and queries.

In this paper, we introduce the Universal Instruction-to-Answer Navigator (Uni-Med) framework to extract instruction-to-answer relationships, thereby making answers both traceable and learnable. Specifically, we design an Instruct-to-Answer Clues Interpreter (IAI), which employs
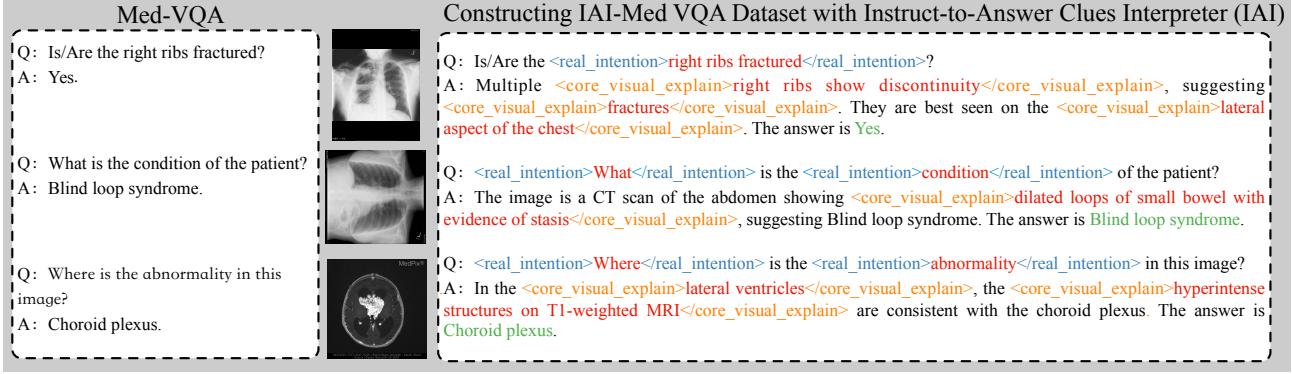
*Figure 1.* Details of the construction of the IAI-Med VQA dataset with the Instruct-to-Answer Clues Interpreter (IAI). We design a Universal-Navigator Prompt (UNP) to guide MLLM to articulate the reasoning behind answers based on the visual content present in medical images and the context provided by existing question-answer pairs. The label "real_intent" is used to label the "real intent" of user instructions and "core_visual_explain" is used to mark the visual clues that support the explanation.

an MLLM to generate visual explanations as reasoning steps. The IAI-Med VQA dataset, created by IAI, enhances the VQA dataset by adding "real intent" labels of instructions and providing corresponding visual explanations, as shown in Figure 1. The "real intent" refers to the key information that users want the model to focus on, reducing interruptions by irrelevant details. To mitigate errors in medical image analysis, we developed the Universal-Navigator Prompt (UNP) within IAI, which guides the MLLM in categorizing answers based on different types, such as questions and organs.

Additionally, Uni-Med integrates a token-level feature enhancement strategy using 'core visual explanation' labels. This approach aligns visual explanations with corresponding image blocks for token-level cut-mix processing, concentrating the model on task-relevant visuals and helping the LLM locate the source of the answer. Furthermore, we designed an Intention-guided Attention (IGA) mechanism that adaptively reduces the attention score for non-core instructions, thereby sharpening the LLM's focus on 'real intent' content to minimize modal interference. The whole process provides direction to the LLM in answering questions, just like a navigator. The proposed method achieves state-of-the-art (SOTA) performance on the SLAKE dataset, exceeding (Zhang et al., 2023b) by 1.22% and 0.92% in closed and overall accuracy, respectively.

The main contributions of this paper are outlined as follows:

- We introduce a Universal Instruction-to-Answer Navigator Learning Framework (Uni-Med) for extracting any instruction-to-answer interaction relationship, which make the answer 'traceable' and 'learnable'.

- We design an Instruct-to-Answer Clues Interpreter (IAI) to generate the IAI-Med VQA dataset, which marks the "real intent" of instructions and generates corresponding visual explanations. To minimize errors

in medical image analysis, we develop an Universal-Navigator Prompt (UNP) to enhance medical image understanding and reasoning of MLLM.

- We implement a task-guided Token-level Cut-Mix (TC-Mix) strategy that leverages visual explanation aligned with user instructions, mapping them to the most relevant blocks in medical images for token level feature enhancement.

## 2. Related Work

**Biomedical Visual Question Answering.** Current Med-VQA approaches typically handle inquiries through classification tasks, sourcing responses from a predetermined set (Binh D. Nguyen, 2019)(Do et al., 2021)(Gong et al., 2022b). While this method performs well for closed-ended questions, it is less effective for clinical open-ended inquiries (Chen et al., 2022; 2023; Yuan et al., 2023). Notably, medical chatbot like LLaVA-Med (Li et al., 2023a) require extensive fine-tuning on large instruction datasets to effectively follow user instructions (Liu et al., 2023). Besides, PMC-VQA (Zhang et al., 2023b) has developed a substantial medical VQA dataset to enhance medical visual comprehension. Nonetheless,exsiting works often treat text and image understanding equally, overlooking the challenge of modal interference. Some exploratory methods (Tascon-Morales et al., 2023; Han et al., 2020) enhance the interpretability of answers by focusing on image location related to the question, but do not filter irrelevant information in the question and may be interfered by it (Liu et al., 2024a). Complex instructions can further prevent the model from focusing on the user's intentions, resulting in undesired outputs (Wei et al., 2022). Recently, (Zhang et al., 2023a) introduced a post-attention method to focus the crucial elements of instructions during training, but it relies on manual marking, and becomes inefficient for large numbers of queries.

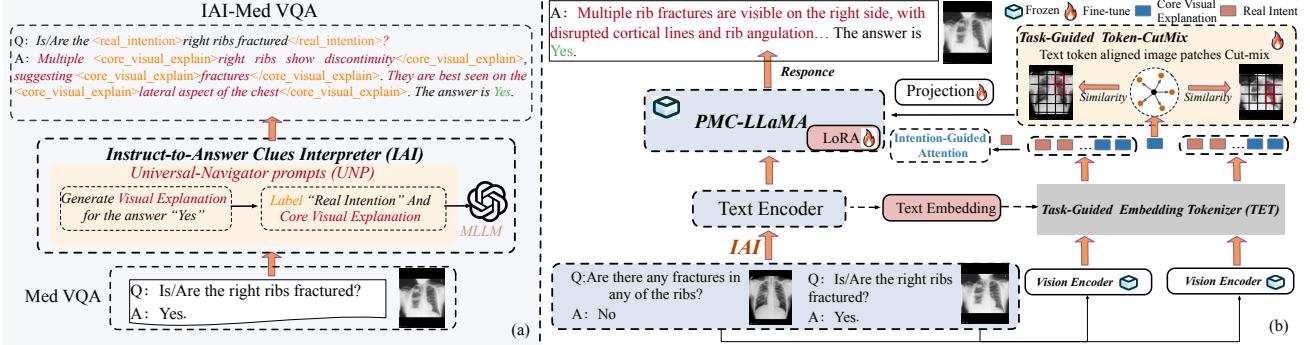**Medical Data Augmentation.** Adapting general VQA mod-

*Figure 2.* The Uni-Med Training Paradigm. (a): The IAI module, where UNP prompts MLLMs to identify instruction's "real intent" and generate visual explanations. (b): TET is a tokenizer that distinguish the text embedding corresponding to "core visual explanation" and "real intent". The most aligned patches are selected to perform feature-level enhancement by TC-Mix. Intention-Guided Attention then focuses LLM (*e.g.*, PMC-LLAMA) on the "real intent" to minimize modal interference.

els to medical applications often leads to overfitting due to data variance and scarcity(Su et al., 2022a; Cao et al., 2022). Previous works (Ray et al., 2019; Tang et al., 2020a; Gupta et al., 2021) have improved model consistency using logically consistent QA pairs, albeit dependent on external dataset. MixUp techniques have been applied in both image and text classification to bolster generalization in downstream tasks (Verma et al., 2019; Guo et al., 2019; Guo, 2020). The introduction of conditional VQAMix (Gong et al., 2022a; 2021) addresses data scarcity by increasing training data diversity (Su et al., 2021d). Simplification of the augmentation process in some methods has led to the generation of new samples, thus contributing to the robustness and diversity of medical VQA models (Tang et al., 2020b; Agarwal et al., 2020; Liu et al., 2024b). However, the (Zhang et al., 2020) approach, despite utilizing Grad-CAM guided, feature-level CutMix(Yun et al., 2019), neglects task-specific instructions.

## 3. Problem Formulation

Med-VQA involves responding to natural language queries about medical visual content, typically sourced from medical imaging modalities such as X-ray, CT, MRI, or microscopy. In the conventional Med-VQA task, the primary goal is to generate a specific answer $\hat{a}_i$ corresponding to a given image $I_i$ and query $q_i$. This objective can be formulated as follows:

$$\hat{a}_i = \Phi_{\text{MedVQA}}(I_i, q_i; \Theta), \quad (1)$$

where $\Phi_{\text{MedVQA}}$ denotes the function that models the answer generation process and $\Theta$ represents the model parameters.

To improve the interpretability of the Med-VQA task, we developed an Instruct-to-Answer Clues Interpreter (IAI) that integrates visual explanations with textual answers. This integration led to the creation of a dataset named IAI-Med VQA, which is defined mathematically in Eq.(2):

$$(\hat{v}_i, \hat{a}_i) = \Phi_{\text{IAI-MedVQA}}(I_i, q_i; \Theta), \quad (2)$$

where $\hat{v}_i$ represents the visual explanation corresponding to $\hat{a}_i$. Details will be described in Section 4.1.

For Med-VQA task, the loss function for textual answers is formulated as:

$$L_{\text{txt}}(\Theta) = -\sum_{t=1}^{T} \log p(\hat{a}_i^t \mid I_i, q_i^{1:T}, \hat{a}_i^{1:t-1}; \Theta), \quad (3)$$

where $T$ is the length of the answer sequence. To enhance interpretability, we introduce a loss function for visual explanations:

$$L_{\text{vis}}(\Theta) = \sum_{j=1}^{J} \|\hat{v}_i^j - v_i^{*j}\|_2^2, \quad (4)$$

where $J$ is the number of elements in the visual explanation and $v_i^{*j}$ is the ground-truth visual explanation for the $j$-th element. The overall training objective combines these components:

$$\Theta^* = \arg\min_{\Theta} \left( L_{\text{txt}}(\Theta) + \lambda L_{\text{vis}}(\Theta) \right), \quad (5)$$

with $\lambda$ as a hyperparameter to balance the two losses. By optimizing Equation (5), our model not only answers medical queries but also provides text-based visual explanations, enhancing the interpretability of answers in Med-VQA.

**Objective.** Our goal is to generate a natural language response $\hat{y}_i$ that effectively encapsulates both the visual content $V_i$ and the task-specific aspects of the query $Q_i$. This is accomplished using our setting, denoted as Uni-Med VQA.

$$\begin{aligned} \hat{y}_i &= \Gamma_{\text{Uni-Med VQA}}(V_i, Q_i^{Task}; \Lambda) \\ &= \Gamma_{\text{dec}}(\Gamma_{\text{img}}(V_i; \lambda_{\text{img}}), \Gamma_{\text{text}}(Q_i^{Task}; \lambda_{\text{text}}); \lambda_{\text{dec}}) \end{aligned} \quad (6)$$

Here, $\Lambda$ represents parameter set of the enhanced VQA model, including $\lambda_{\text{img}}, \lambda_{\text{text}}, \lambda_{\text{dec}}$. The model integrates the visual and query processing modules to output contextually relevant and task-aligned answers.

# 4. Method

The Uni-Med architecture, depicted in Figure 2, begins with utilizing the IAI module to construct the IAI-Med VQA dataset. In this dataset, each answer is accompanied by a visual explanation, with the essential parts of both instructions and explanations highlighted. Following this, we capture embeddings based on these labelled explanations and instructions. These embeddings are then employed for two key purposes: token-level feature enhancement, as detailed in Section 4.2 and assisting the LLM in focusing on the core aspects of the instruction, elaborated in Section 4.3.

## 4.1. Instruct-to-Answer Clues Interpreter

In order to provide evidence-based answers, we designed an IAI to generate an IAI-Med VQA dataset, which marked the "real intent" of instructions and generated corresponding visual explanations.

**Visual Explanation Generation.** For a given image $I_i$ and question-answer pair $(Q_i, A_i)$, we generate a visual explanation $V_i$ using MLLM without parameter updates:

$$V_{\text{cot}_i} = \text{MLLM}(I_i, Q_i, A_i; \text{UNP}), \tag{7}$$

where UNP represents Universal-Navigator prompts, as illustrated in Figure 4, which instruct the MLLM to generate Chain of Thought (COT) explanations. $T_{\text{cot}_i}$ represents the textual output that combines the visual explanation and reasoning for the answer.

To establish the relationship between answers and vision and to guide the MLLM to focus on the 'real intent', it is essential to annotate the core instructions using UNP (refer to Appendix Figure 12 for comprehensive details).

**Core Instruction and Explanation Labeling.** We define the process of labeling the real intent of instructions and associated visual explanations as:

$$(q_i^R, V_i^R) = \text{MLLM}(I_i, Q_i, V_{\text{cot}_i} + A_i; \text{UNP}), \tag{8}$$

where $q_i^R$ denotes the question annotated with 'real_intent' labels and $V_i^R$ includes 'core_visual_explain' labels for visual clues in the image that support the answer.

**Generation of IAI-Med VQA Dataset.** Following the acquisition of the results from IAI, the structure of our optimized Med-VQA task is delineated as Eq.(9):

$$\begin{aligned} \text{Question:} \quad & q_i^R < \text{STOP} > \\ \text{Answer:} \quad & V_{\text{cot}_i} + a_i^R < \text{STOP} > \end{aligned} \tag{9}$$

In this way, we developed a dataset with traceable instruction-to-answer paths to enhance the interpretability of Med-VQA. It augments the VQA-RAD dataset with 'real intent' annotations and visual explanations. To advance Med-VQA research, we have made the IAI-Med VQA dataset publicly available.
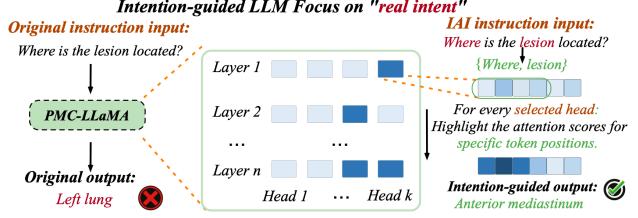


*Figure 3.* The details of Intention-guided Attention.

## 4.2. Task-guided Token-Level Cut-Mix

Building on the IAI module, we introduce the TC-Mix strategy (detailed in Algorithm 1), to further enhance feature-level understanding. As shown in Figure 7, when differentiating between "pneumonia" and "pleural effusion", the "real intent" is to identify core visual clues such as "homogenous opacity" (marked in a red box). TC-Mix improves the model's accuracy in identifying the core visual clues necessary for accurate diagnosis by token-level cut mixing. This module dynamically aligns core visual explanations $A_j^M$ with corresponding image patches, enabling Uni-Med to identify finer-grained differences and ensure each response is 'traceable', as shown in Figure 7.

**Image Group and decomposition.** In our approach, medical images are categorized by organ type $O$ and problem category $P$, resulting in distinct groups $\mathcal{G} = \{G_{o,p}\}$ (as in Appendix Figure 15). Each group $G_{o,p}$ is associated with a set of marked answers $\mathcal{A}_{o,p} = \{A_i^{o,p}\}_{i=1}^{N_{o,p}}$, where $N_{o,p}$ represents the number of question-answer pairs in the group and each answer $A_i^{o,p}$ contains M answer-labeled text tokens $A_j^{M_{o,p}}$. Here, $j$ is an index representing the position of a specific answer within the marked answers $\mathcal{A}_i^{o,p}$ and M represents the number of marked tokens in $\mathcal{A}_i^{o,p}$.

Firstly, each medical image $I_i^{o,p}$ is decomposed into fixed-size patches $P_{I_i}^{o,p}$ of a fixed size, as detailed in Eq.(10):

$$\begin{aligned} P_{I_i}^{o,p} = I_i^{o,p}[ih : (i+1)H, jW : (j+1)W] \mid \\ i \in \{0, 1, \ldots, M-1\}, j \in \{0, 1, \ldots, N-1\}. \end{aligned} \tag{10}$$

We denote a patch by $p_{i,k} \in \mathbb{R}^{b_h \times b_w}$, where $k$ is the index of the patch and $b_h$ and $b_w$ are the predetermined height and width of each patch, respectively. The total number of patches $K$ in image $I_i^{o,p}$ is determined by the ratio of the image dimensions $H \times W$ to the patch dimensions $b_h \times b_w$.

$$P_{I_i}^{o,p} = \{p_{i,k} | k \in [1, K]\}, \tag{11}$$

where $P_{I_i}^{o,p}$ represents the set of patches into which image $I_i^{o,p}$ is decomposed.

**Semantic Similarity-Based Patch Ranking.** As shown in Figure 5, for each image patch $p_{i,k}$, we compute its feature vector $v_{i,k}$ using an image encoder $\mathcal{E}_{\text{img}}$. And for each text token, we compute $t_j$ using a text encoder $\mathcal{E}_{\text{txt}}$:

$$v_{i,k} = \mathcal{E}_{\text{img}}(p_{i,k}), \quad t_j = \mathcal{E}_{\text{txt}}(A_j^M), \tag{12}$$

Universal-Navigator prompt in IAI to generate Visual Explanation based on current conversation and image clues.

*payload* = {
    "model": "gpt-4-vision-preview", "messages": [{"role": "system","content": (
        """You are now operating as a medical imaging specialist. Your role is to analyze medical images and provide concise and precise explanations for specific findings... Remember to:
        1. Clearly identify the modality of the medical image (e.g., X-Ray, MRI).
        2. Pinpoint the location and size of any notable features or abnormalities...
        4. Refrain from overinterpreting arrows or labels; focus instead on the pathology they may indicate...
        6. Avoid ambiguous language and ensure your response is decisive and clinically relevant...
    The answer is the real label. You need to analyze and keep the basic facts consistent with the answer.
    Answers based on COT should have different focus points for different question types, such as:
    -For MODALITY type questions, answers should focus on confirming what type of medical imaging (e.g., X-ray, MRI, CT scan, etc.).
    -For ORGAN type questions, answers should focus on the organs and their abnormalities visible in the image···"""
        )},{ "role": "user","content": [{"type": "text","text": ( Here is a new medical image along with a question.... f"The current question and answer of uploaded image is {question} and {answer}, The question type is {question_type} and the organ type is {image_organ}. You can refer to the reply methods in the template:{template}")},{"type": "image_url","image_url": {"url": f"data:image/jpeg;base64,{base64_image}"}}]}],"max_tokens": 200}

*Figure 4.* The simplified version of prompt MLLM to generate task related explanation.

where $A_j^M$ represents the marked text token relevant to image patch $p_{i,k}$. Then we using projection heads to project them to a common feature space:

$$\tilde{v}_{i,k} = f_{img}(v_{i,k}), \quad \tilde{t}_j = f_{txt}(t_j), \tag{13}$$

yielding normalized feature vectors for comparison.

The cosine similarity between the patch $p_{i,k}$ and text token $A_j^M$ is computed as:

$$s_{ij} = \tilde{v}_{i,k}^\top \cdot \tilde{t}_j, \tag{14}$$

The cosine similarity is then normalized across all tokens and patches using a temperature-scaled softmax function to get semantic similarity scores. The semantic similarity from patch to text token is given by:

$$\sigma_{ij}^{p2t} = \frac{\exp(s_{ij}/\tau)}{\sum_{k=1}^{N_{patches}} \exp(s_{ik}/\tau)}, \tag{15}$$

where $\tau$ is the temperature parameter. The final semantic score for a patch-token pair is calculated as:

$$S_{i,k}^{'} = \lambda_1 \cdot \sigma_{ij}^{p2t} + \lambda_2 \cdot \sigma_{ji}^{t2p}, \tag{16}$$

where $\sigma_{ji}^{t2p}$ is the reverse semantic similarity from text token to patch and $\lambda_1$ and $\lambda_2$ are weights to balance the contributions of patch-to-text and text-to-patch similarities, respectively.

**Group Based Token-CutMix.** For each image pair within group $G_{o,p}$, we iterate the selection of top-N similar patches and the generation of binary masks $m$. The binary mask $m_{i,k}$ for each patch $k$ in image $I_i^{o,p}$ is defined using the top-N similarity scores $S_{i,k}^{'}$:
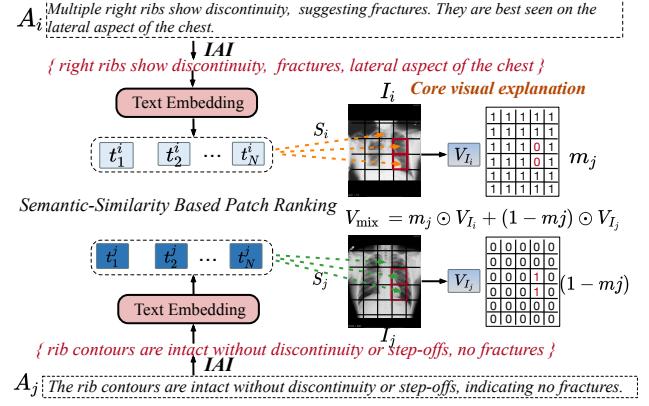


*Figure 5.* The details of Task-guided Token-Level Cut-Mix.

$$m_j = \begin{cases} 0 & \text{if } j \in \{p_i, k\} \& S_{i,k}^{'} \in \left\{ \text{top}_N(S_{i,k}^{'}) \right\}, \\ 1 & \text{otherwise.} \end{cases} \tag{17}$$

In Eq.(17), $j \in [1, H \times W]$ indexes the spatial dimensions of mask $m_j$. The TC-Mix operation is formulated as follows:

$$V_{mix} = m_j \odot V_{I_i} + (1 - m_j) \odot V_{I_j}, \tag{18}$$

where $\odot$ represents element-wise multiplication and $V_{I_i}$, $V_{I_j}$ are the feature vectors from respective images in the group $\mathcal{G}$.

### 4.3. Intention-guided LLM Focus on "real intent"

To reduce modal interference and further improve the model's ability to follow instructions, we implement an intention-guided attention mechanism, directing the model's
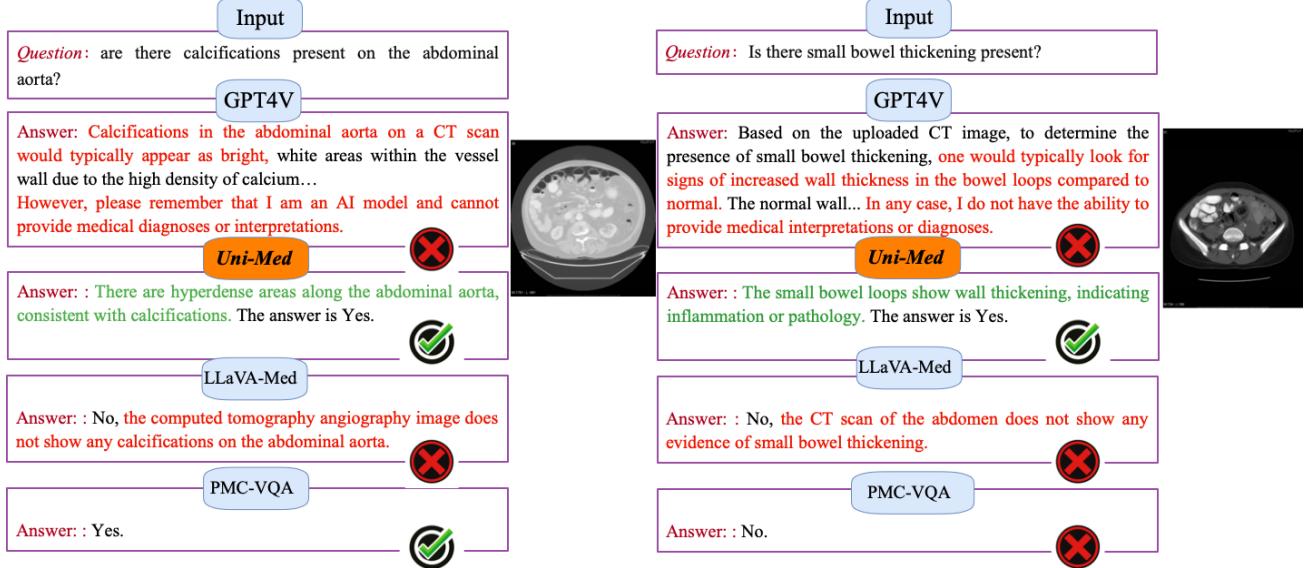
*Figure 6.* Qualitative analysis between different MLLMs and Our Uni-Med.

focus towards 'real intention' parts $q_i^{TR}$ of the instructions, detailed in Figure 3.The "real intention" of a question often guides us to find corresponding visual clues in an image. TC-Mix improves the model's accuracy in identifying the core visual clues necessary for accurate diagnosis by token-level cut mixing. This approach adjusts attention scores $A^{(l,k)}$ within each head $k$ of layer $l$ in the multi-head attention (MHA) layers of the LLM like PMC-LLaMA.

**Intention-guided Attention.** We identify the set of indices $S_q$ corresponding to the tokens $q_i^{TR}$ within $q_i^R$ and apply an attention projection $\mathbf{P}$ to adjust the attention scores in the MHA layers:

$$H_{\text{IGA}}^{(l,k)} = \mathbf{P}(A^{(l,k)})V, \qquad (19)$$

where $\mathbf{P}$ is the projection function, as defined in Eq.(20):

$$[\mathbf{P}(A)]_{ij} = \begin{cases} \beta A_{ij}/D_i & \text{if } j \notin S_q \\ A_{ij}/D_i & \text{otherwise} \end{cases} \qquad (20)$$

Here, $0 < \beta < 1$ is a scaling factor to reduce attention scores for tokens outside $q_i^R$. We define $D_i$ for normalization as $D_i = \sum_{j \notin S} \beta A_{ij} + \sum_{j \in S} A_{ij}$ ensuring the scores sum to one. The value matrix $V$ in the MHA mechanism, combines with the attention scores to contribute to the final output of each head.

By reducing the attention scores for tokens not marked as 'real intention', our method maintains the relative importance of emphasized segments, avoids uniform attention distribution that can lead to information loss and ensures numerical stability. This effectively aligns model output with 'real intention' content. As illustrated in Figure 1, for questions like "What is the condition of the patient?" and "Where is the abnormality in this image?", our approach prioritizes focus on critical subsets [What, condition] and

[Where, abnormality], directing the LLM's focus towards the real intent of the instructions, such as 'condition' and 'abnormality'. This minimizes distractions from non-essential details, ensuring the model's attention remains on pertinent instruction elements. By employing IGA, we direct the model's attention to the important details of the instructions, mitigating the issue of LLMs struggling to capture key information amidst length and complex instructions.

## 5. Experiment Results

### 5.1. Experiment Settings

**Dataset.** We use the PMC-VQA dataset (Zhang et al., 2023b), which includes 227K VQA pairs from 149K images, adhered to the experimental configurations as described in (Nguyen et al., 2019a). For fine-tuning, we used two medical datasets: VQA-RAD (Nguyen et al., 2019a), consisting of 314 radiology images and 3,064 clinician-curated question-and-answer pairs; and SLAKE (Liu et al., 2021b), which offers 642 radiology images and 14K question-and-answer samples, of which we used 70% for training and 30% for testing. For these two datasets, the experimental setup in (Zhang et al., 2023b) were followed.

**Implementation details.** In our implementation, the PMC-CLIP visual backbone (Lin et al., 2023) is kept frozen throughout the training process. We align our model's optimization strategy, learning rate, number of training epochs, weight decay, loss weights and other hyperparameters with established methods to ensure a fair comparison. The model is trained using the AdamW optimizer, combined with a cosine learning rate scheduler, across 8 Tesla V100 GPUs over 8,000 steps. We set a global batch size of 128 and a peak learning rate of 2e-5 to optimize performance.

**Algorithm 1** Token-Level Cut-Mix (TC-Mix)

1: **Input:** Image groups $\mathcal{G} = \{G_{o,p}\}$, Real Intent Marked answers $\mathcal{A}_{o,p} = \{A_i^{o,p}\}_{i=1}^{N_{o,p}}$
2: **Output:** Cut-Mixed feature vectors $V_{\text{mix}}$
3: **for** each group $G_{o,p}$ in $\mathcal{G}$ **do**
4:     **for** each answer $A_i^{o,p}$ in $\mathcal{A}_{o,p}$ **do**
5:         Concatenate marked answer tokens $A_j^{M_{o,p}} \leftarrow$ Concat$(a_1^{o,p}, a_2^{o,p}, ..., a_n^{o,p})$ from $A_i^{o,p}$
6:         Compute concatenated text embedding $t_j^{o,p} \leftarrow \mathcal{E}_{\text{txt}}(A_j^{M_{o,p}})$
7:     **end for**
8:     **for** each image $I_i^{o,p}$ in $G_{o,p}$ **do**
9:         Decompose $I_i^{o,p}$ into K patches $P_{I_i}^{o,p} = \{p_{i,k}\}_{k=1}^{K}$
10:        Extract image features $v_{i,k} \leftarrow \mathcal{E}_{\text{img}}(p_{i,k})$, for all $k$
11:        Calculate semantic similarity scores $S'_{i,k} \leftarrow \kappa(v_{i,k}, t_j^{o,p})$, for all $k$
12:        Select top-N patches $T_S^N \leftarrow \text{top}_N(S'_{i,k})$
13:     **end for**
14:     **for** each pair $(I_i, I_j)$ in $G_{o,p}$ **do**
15:        Generate binary masks $m_i, m_j$ for $T_S^N$
16:        Perform Cut-Mix: $V_{\text{mix}} \leftarrow (m_i \odot V_{I_i}) + (1 - m_j) \odot V_{I_j}$
17:     **end for**
18: **end for**
19: **Return:** $V_{\text{mix}}$

*Q:Is the trachea midline?*
*A:<core_visual_explanation>The trachea<core_visual_explanation> appears <core_visual_explanation>straight and centered between the clavicles</core_visual_explanation>, indicating it is midline.The answer is yes.*



*Q:Is this pneumonia vs. pleural effusion?*
*A:There is <core_visual_explanation>homogenous opacity</core_visual_explanation> in the <core_visual_explanation>right hemithorax</core_visual_explanation> with visible <core_visual_explanation>loss of the right diaphragmatic silhouette</core_visual_explanation>, indicating pleural effusion.*



*Figure 7.* Visualization result of the TC-Mix On VQA-RAD.

*Table 5.* Stability and GPU throughput comparison.

| Method | Closed | Open | Overall | GPU Throughput (samples/sec) |
|---|---|---|---|---|
| **VQA-RAD** | | | | |
| LLaVA | 65.02 ± 0.12 | 49.95 ± 0.18 | - | 2.86 |
| LLaVA-Med | 84.15 ± 0.13 | 61.5 ± 0.16 | - | 3.23 |
| MedVInT-PMC-VQA | 86.75 ± 0.1 | 73.65 ± 0.15 | 81.55 ± 0.12 | 2.57 |
| Uni-Med | 87.15 ± 0.08 | 74.18 ± 0.12 | 81.95 ± 0.09 | 2.48 |
| **SLAKE** | | | | |
| LLaVA | 63.2 ± 0.14 | 78.1 ± 0.17 | - | 2.78 |
| LLaVA-Med | 85.3 ± 0.12 | 83.05 ± 0.15 | - | 3.12 |
| MedVInT-PMC-VQA | 86.25 ± 0.1 | 84.45 ± 0.13 | 85.15 ± 0.11 | 2.51 |
| Uni-Med | 87.48 ± 0.07 | 85.25 ± 0.11 | 86.05 ± 0.10 | 2.42 |

dataset, IAI-Med facilitates an increase of over 5% for both question formats across different model sizes. These findings highlight the effectiveness of IAI-Med in improving MLLMs' interpretative and diagnostic capabilities.

**Effects of Visual Interpretation Position.** According to Table 2, the positioning of visual interpretations significantly influences model performance. On the VQA-RAD dataset, situating visual interpretations before the answer (IAI-Med-pre) enhances accuracy for closed questions by 0.95%, while placing them after the answer (IAI-Med) leads to a 1.82% improvement. For open-ended questions, IAI-Med-pre offers a slight advantage over IAI-Med, indicating that while IAI-Med adheres to a traditional, reasoning-oriented chain of thought, IAI-Med-pre might facilitate a more comprehensive understanding.

**Impact of UNP in IAI module.** To examine the influence of UNP on IAI, we performed experiments as shown in Table 3, with the baseline being MedVInT-PMC-VQA. On the VQA-RAD dataset, the removal of UNP caused a 1.18% reduction in overall accuracy. Similarly, on the SLAKE dataset, the exclusion of UNP led to a 1.57% decrease in overall accuracy. These observations underscore the critical contribution of UNP to the IAI.

**Effects of IAI, IGA and TC-Mix.** Our ablation study (Table 4) reveals the significance of each component in the framework. Analysis indicates that excluding the IAI mod-

*Table 4.* Effects of IAI, IGA and TC-Mix.

| Method | VQA-RAD | | | SLAKE | | |
|---|---|---|---|---|---|---|
| | Closed | Open | Overall | Closed | Open | Overall |
| baseline | 86.8 | 73.7 | 81.6 | 86.3 | 84.5 | 85.2 |
| w/o IAI | 87.12 | 73.85 | 81.78 | 87.40 | 84.78 | 85.82 |
| w/o TC-Mix | 87.02 | 74.05 | 81.86 | 87.28 | 85.06 | 85.77 |
| w/o IGA | 87.11 | 73.92 | 81.72 | 87.42 | 84.60 | 85.73 |
| Full | **87.22** | **74.21** | **82.05** | **87.52** | **85.34** | **86.12** |

trachea, straight and centered between the clavicles" with their respective image patches, thereby focusing the model's attention on visuals pertinent to the task. This method involves selectively employing TC-Mix on image patches that coincide with these visual explanations, thereby sharpening the model's focus on essential visuals. Such precise alignment, facilitated by TC-Mix processing, significantly boosts the model's ability to adapt to unfamiliar content while maintaining the clarity of its responses.

## 5.5. Ablation Studies

**Impact of using IAI-Med VQA Dataset.** As illustrated in Table 2, using IAI-Med VQA dataset enhances the performance of general MLLMs like LLaVA-1.5. For the 7B model, using IAI-Med improves accuracy by 2.71% for open and 1.82% for closed questions on the VQA-RAD dataset. The 13B model sees even larger gains of 3.11% and 2.33% for open and closed questions, respectively. On the SLAKE
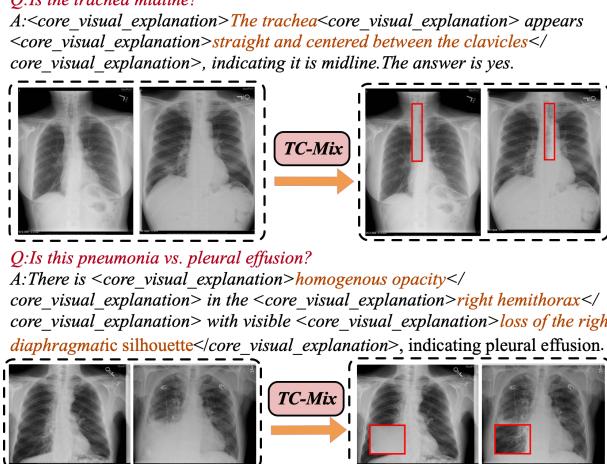
ule does not hinder performance gains over the baseline, due to the complementary support of other modules. TC-Mix, which enhances the model's ability to focus on diverse features, impacts performance more than IGA . The combined inclusion of IAI, IGA, and TC-Mix markedly improves outcomes across various question types, underscoring their collective importance in boosting the model's effectiveness.

**Stability and GPU throughput comparison.** In our study, we performed ten experiments to evaluate Uni-Med against other SOTA methods, focusing on stability and GPU throughput. The results, detailed in Table 5, reveal that Uni-Med consistently exhibits a standard deviation within the 1% range for both VQA-RAD and SLAKE datasets, indicating higher reproducibility and reliability compared to its counterparts. Moreover, Uni-Med showcases competitive GPU throughput, highlighting its efficiency alongside its performance robustness.

## 6. Conlusion

In this study, we introduce the Uni-Med framework, an approach that significantly enhances the interpretation of complex medical images through user instructions. Utilizing the IAI for identifying the 'real intent' behind queries and generating precise visual explanations, alongside the Token-Level Cut-Mix module for ensuring traceable and learnable answers, Uni-Med marks a significant advance in Med-VQA research. Extensive experiments on public datasets and diverse settings demonstrate its effectiveness. By making the IAI-Med VQA dataset publicly available, we aim to foster future advancements in interpreting complex medical images, thereby contributing to the broader field of medical diagnostics and research.

## Acknowledgements

## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. All the work and data are based on existing public datasets and methods, thus there are no potential adverse societal consequences.

## References

Agarwal, V., Shetty, R., and Fritz, M. Towards causal vqa: Revealing and reducing spurious correlations by invariant and covariant semantic editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

Binh D. Nguyen, Thanh-Toan Do, B. X. N. T. D. E. T. Q. D. T. Overcoming data limitation in medical visual question answering. In *MICCAI*, 2019.

Cao, Y., Su, X., Tang, Q., You, S., Lu, X., and Xu, C. Searching for better spatio-temporal alignment in few-shot action recognition. *Advances in Neural Information Processing Systems*, 35:21429–21441, 2022.

Cao, Y., Tang, Q., Yang, F., Su, X., You, S., Lu, X., and Xu, C. Re-mine, learn and reason: Exploring the cross-modal semantic correlations for language-guided hoi detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 23492–23503, 2023.

Cao, Y., Tang, Q., Su, X., Chen, S., You, S., Lu, X., and Xu, C. Detecting any human-object interaction relationship: Universal hoi detector with spatial prompt learning on foundation models. *Advances in Neural Information Processing Systems*, 36, 2024.

Chen, J., Yang, D., Jiang, Y., Lei, Y., and Zhang, L. Miss: A generative pretraining and finetuning approach for med-vqa, 2024.

Chen, Z., Du, Y., Hu, J., Liu, Y., Li, G., Wan, X., and Chang, T.-H. Multi-modal masked autoencoders for medical vision-and-language pre-training. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2022.

Chen, Z., Diao, S., Wang, B., Li, G., and Wan, X. Towards unifying medical vision-and-language pre-training via soft prompts. *arXiv preprint arXiv:2302.08958*, 2023.

Cong, F., Xu, S., Guo, L., and Tian, Y. Caption-aware medical vqa via semantic focusing and progressive cross-modality comprehension. In *Proceedings of the 30th ACM International Conference on Multimedia*, pp. 3569–3577, 2022.

Do, T., Nguyen, B. X., Tjiputra, E., Tran, M., Tran, Q. D., and Nguyen, A. Multiple meta-model quantifying for medical visual question answering. In *MICCAI*, 2021.

Eslami, S., de Melo, G., and Meinel, C. Does clip benefit visual question answering in the medical domain as much as it does in the general domain? *arXiv preprint arXiv:2112.13906*, 2021.

Gong, H., Chen, G., Liu, S., Yu, Y., and Li, G. Cross-modal self-attention with multi-task pre-training for medical visual question answering. pp. 456–460. ACM, 2021. doi: 10.1145/3460426.3463584. URL https://doi.org/10.1145/3460426.3463584.

Gong, H., Chen, G., Mao, M., Li, Z., and Li, G. Vqamix: Conditional triplet mixup for medical visual question answering. *IEEE Trans. Medical Imaging*, 41(11):3332–3343, 2022a. doi: 10.1109/TMI.2022.3185008. URL https://doi.org/10.1109/TMI.2022.3185008.

Gong, H., Chen, G., Mao, M., Li, Z., and Li, G. Vqamix: Conditional triplet mixup for medical visual question answering. *IEEE Trans. on Medical Imaging*, 2022b.

Guo, H. Nonlinear mixup: Out-of-manifold data augmentation for text classification. pp. 4044–4051. AAAI Press, 2020. doi: 10.1609/AAAI.V34I04.5822. URL https://doi.org/10.1609/aaai.v34i04.5822.

Guo, H., Mao, Y., and Zhang, R. Augmenting data with mixup for sentence classification: An empirical study. *CoRR*, abs/1905.08941, 2019. URL http://arxiv.org/abs/1905.08941.

Gupta, D., Suman, S., and Ekbal, A. Hierarchical deep multimodal network for medical visual question answering. *Expert Systems with Applications*, 164:113993, 2021.

Han, T., Adams, L. C., Papaioannou, J.-M., Grundmann, P., Oberhauser, T., Löser, A., Truhn, D., and Bressem, K. K. Medalpaca – an open-source collection of medical conversational ai models and training data, 2023.

Han, W., Huang, H., and Han, T. Finding the evidence: Localization-aware answer prediction for text visual question answering. In *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020*, pp. 3118–3131, 2020.

Kafle, K., Yousefhussien, M., and Kanan, C. Data augmentation for visual question answering. In *Proceedings of the 10th International Conference on Natural Language Generation*, pp. 198–202, 2017.

Khare, Y., Bagal, V., Mathew, M., Devi, A., Priyakumar, U. D., and Jawahar, C. Mmbert: Multimodal bert pretraining for improved medical vqa. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pp. 1033–1036. IEEE, 2021.

Li, C., Wong, C., Zhang, S., Usuyama, N., Liu, H., Yang, J., Naumann, T., Poon, H., and Gao, J. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. *arXiv preprint arXiv:2306.00890*, 2023a.

Li, P., Liu, G., He, J., Zhao, Z., and Zhong, S. Masked vision and language pre-training with unimodal and multimodal contrastive losses for medical visual question answering. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 374–383. Springer, 2023b.

Li, W., Su, X., You, S., Wang, F., Qian, C., and Xu, C. Diffnas: Bootstrapping diffusion models by prompting for better architectures. In *2023 IEEE International Conference on Data Mining (ICDM)*, pp. 1121–1126. IEEE, 2023c.

Li, Y., Li, Z., Zhang, K., Dan, R., Jiang, S., and Zhang, Y. Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge, 2023d.

Lin, W., Zhao, Z., Zhang, X., Wu, C., Zhang, Y., Wang, Y., and Xie, W. Pmc-clip: Contrastive language-image pre-training using biomedical documents. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2023. URL https://api.semanticscholar.org/CorpusID:257496659.

Liu, B., Zhan, L.-M., and Wu, X.-M. Contrastive pre-training and representation distillation for medical visual question answering based on radiology images. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part II 24*, pp. 210–220. Springer, 2021a.

Liu, B., Zhan, L.-M., Xu, L., Ma, L., Yang, Y., and Wu, X.-M. Slake: A semantically-labeled knowledge-enhanced dataset for medical visual question answering, 2021b.

Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning, 2023.

Liu, N. F., Lin, K., Hewitt, J., Paranjape, A., Bevilacqua, M., Petroni, F., and Liang, P. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173, 2024a.

Liu, P., Ji, L., Zhang, X., and Ye, F. Pseudo-bag mixup augmentation for multiple instance learning-based whole slide image classification. *IEEE Transactions on Medical Imaging*, pp. 1–1, 2024b. doi: 10.1109/TMI.2024.3351213.

Moor, M., Huang, Q., Wu, S., Yasunaga, M., Dalmia, Y., Leskovec, J., Zakka, C., Reis, E. P., and Rajpurkar, P. Med-flamingo: a multimodal medical few-shot learner. In *Machine Learning for Health (ML4H)*, pp. 353–367. PMLR, 2023.

Nguyen, B. D., Do, T.-T., Nguyen, B. X., Do, T., Tjiputra, E., and Tran, Q. D. Overcoming data limitation in medical visual question answering, 2019a.

Nguyen, B. D., Do, T.-T., Nguyen, B. X., Do, T., Tjiputra, E., and Tran, Q. D. Overcoming data limitation in medical visual question answering. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part IV 22*, pp. 522–530. Springer, 2019b.

OpenAI. Gpt-4v(ision) system card, 2023. URL https://cdn.openai.com/papers/GPTV_System_Card.pdf.

Ray, A., Sikka, K., Divakaran, A., Lee, S., and Burachas, G. Sunny and dark outside?! improving answer consistency in VQA through entailed question generation. *CoRR*, abs/1909.04696, 2019. URL http://arxiv.org/abs/1909.04696.

Su, X., Huang, T., Li, Y., You, S., Wang, F., Qian, C., Zhang, C., and Xu, C. Prioritized architecture sampling with monto-carlo tree search. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10968–10977, 2021a.

Su, X., You, S., Huang, T., Wang, F., Qian, C., Zhang, C., and Xu, C. Locally free weight sharing for network width search. *arXiv preprint arXiv:2102.05258*, 2021b.

Su, X., You, S., Wang, F., Qian, C., Zhang, C., and Xu, C. Bcnet: Searching for network width with bilaterally coupled network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2175–2184, 2021c.

Su, X., You, S., Zheng, M., Wang, F., Qian, C., Zhang, C., and Xu, C. K-shot nas: Learnable weight-sharing for nas with k-shot supernets. In *International Conference on Machine Learning*, pp. 9880–9890. PMLR, 2021d.

Su, X., You, S., Xie, J., Wang, F., Qian, C., Zhang, C., and Xu, C. Searching for network width with bilaterally coupled network. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022a.

Su, X., You, S., Xie, J., Zheng, M., Wang, F., Qian, C., Zhang, C., Wang, X., and Xu, C. Vitas: Vision transformer architecture search. In *European Conference on Computer Vision*, pp. 139–157. Springer, 2022b.

Tang, R., Ma, C., Zhang, W. E., Wu, Q., and Yang, X. Semantic equivalent adversarial data augmentation for visual question answering. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIX 16*, pp. 437–453. Springer, 2020a.

Tang, R., Ma, C., Zhang, W. E., Wu, Q., and Yang, X. Semantic equivalent adversarial data augmentation for visual question answering. In *European Conference on Computer Vision (ECCV)*, 2020b.

Tascon-Morales, S., Márquez-Neila, P., and Sznitman, R. Localized questions in medical visual question answering. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 361–370. Springer, 2023.

Verma, V., Lamb, A., Beckham, C., Najafi, A., Mitliagkas, I., Lopez-Paz, D., and Bengio, Y. Manifold mixup: Better representations by interpolating hidden states. volume 97 of *Proceedings of Machine Learning Research*, pp. 6438–6447. PMLR, 09–15 Jun 2019. URL https://proceedings.mlr.press/v97/verma19a.html.

Wang, H., Liu, C., Xi, N., Qiang, Z., Zhao, S., Qin, B., and Liu, T. Huatuo: Tuning llama model with chinese medical knowledge, 2023.

Wang, Z., Wu, Z., Agarwal, D., and Sun, J. Medclip: Contrastive learning from unpaired medical images and text. *arXiv preprint arXiv:2210.10163*, 2022.

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., Zhou, D., et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35: 24824–24837, 2022.

Xu, H., Su, X., You, S., Huang, T., Wang, F., Qian, C., Zhang, C., Xu, C., Wang, D., and Sowmya, A. Data agnostic filter gating for efficient deep networks. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 3503–3507. IEEE, 2022.

Yan, Z., Zhang, K., Zhou, R., He, L., Li, X., and Sun, L. Multimodal chatgpt for medical applications: an experimental study of gpt-4v, 2023.

Ye, Q., Xu, H., Ye, J., Yan, M., Hu, A., Liu, H., Qian, Q., Zhang, J., Huang, F., and Zhou, J. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration, 2023.

Yuan, Z., Jin, Q., Tan, C., Zhao, Z., Yuan, H., Huang, F., and Huang, S. Ramm: Retrieval-augmented biomedical visual question answering with multi-modal pre-training. *arXiv preprint arXiv:2303.00534*, 2023.

Yun, S., Han, D., Oh, S. J., Chun, S., Choe, J., and Yoo, Y. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.

Zhan, C., Zhang, Y., Lin, Y., Wang, G., and Wang, H. Unidcp: Unifying multiple medical vision-language tasks via dynamic cross-modal learnable prompts. *arXiv preprint arXiv:2312.11171*, 2023.

Zhan, L.-M., Liu, B., Fan, L., Chen, J., and Wu, X.-M. Medical visual question answering via conditional reasoning. In *Proceedings of the 28th ACM International Conference on Multimedia*, pp. 2345–2354, 2020.

Zhang, Q., Singh, C., Liu, L., Liu, X., Yu, B., Gao, J., and Zhao, T. Tell your model where to attend: Post-hoc attention steering for llms, 2023a.

Zhang, X., Wu, C., Zhao, Z., Lin, W., Zhang, Y., Wang, Y., and Xie, W. Pmc-vqa: Visual instruction tuning for medical visual question answering. *arXiv preprint arXiv:2305.10415*, 2023b.

Zhang, Y., He, B., and Sun, L. Grad-cam guided progressive feature cutmix for classification. *CoRR*, abs/2007.08779, 2020. URL https://arxiv.org/abs/2007.08779.

Universal-Navigator prompt in IAI to generate Visual Explanation based on current conversation and image clues.

```
payload = {
    "model": "gpt-4-vision-preview",
    "messages": [{
        "role": "system","content": (
        """You are now operating as a medical imaging specialist. Your role is to analyze medical images and provide concise and
precise explanations for specific findings.
        Please focus on the visual content of the image and use your medical knowledge to support your reasoning.
        Keep your explanations are directly related to the image's visual aspects.
        Avoid any reliance on text descriptions that are not corroborated by visual evidence in the image.
        Remember to:
        1. Clearly identify the modality of the medical image (e.g., X-Ray, MRI).
        2. Pinpoint the location and size of any notable features or abnormalities.
        3. Distinguish between different types of pathologies or structures based on their visual characteristics.
        4. Refrain from overinterpreting arrows or labels; focus instead on the pathology they may indicate.
        5.Provide a rationale for your answers that demonstrates an understanding of medical imaging conventions.
        6. Avoid ambiguous language and ensure your response is decisive and clinically relevant.
        7. Do not require multi-round clarification; your first response should be as accurate as possible.
    When answering, begin with a direct response to the question, followed by a brief explanation based on the visual evidence.
    Your explanation should be rooted in the image provided and not assume information beyond what is visible.
    The answer is the real label. You need to analyze and keep the basic facts consistent with the answer.
    Answers based on COT should have different focus points for different question types, such as:
    -For MODALITY type questions, answers should focus on confirming what type of medical imaging (e.g., X-ray, MRI, CT scan, etc.).
    -For ORGAN type questions, answers should focus on the organs and their abnormalities visible in the image.
    -For PRES type questions, the focus needs to be on confirming or ruling out the presence of specific pathological phenomena.
    -For POS type questions, the focus is on describing the exact location where the anomaly was found.
    -For ABN type questions, the focus on explaining why a specific pathological abnormality is or is not present in the image…
        """
    )},{
        "role": "user",
        "content": [{"type": "text",
            "text": ( Here is a new medical image along with a question.Please provide an answer and reasoning that strictly pertains
to the visual aspects of this image, without referring to any previous conversation content. Your explanation should directly address the
question and describe the supporting visual evidence found in the image. Be clear, concise, and ensure your response is relevant and
specific to the image at hand.
        f"The current question and answer of uploaded image is {question} and {answer}, The question type is {question_type} and the
organ type is {image_organ}. You can refer to the reply methods in the template:{template}"
        )},{"type": "image_url",
            "image_url": {
                "url": f"data:image/jpeg;base64,{base64_image}"}}]}],
    "max_tokens": 200}
```

*Figure 8.* Universal-Navigator prompt in IAI to generate Visual Explanation based on current conversation and image clues.

## A. *Appendix.*

### A.1. Using Universal-Navigator prompt to generate Visual Explanation

In this section, we design an Instruct-to-Answer Clues Interpreter (IAI), which employs an MLLM (*i.e.*, GPT4V) to generate visual explanations as reasoning steps, as shown in Figure 8. We require GPT4V to act as a medical image expert to accurately analyze and interpret medical images based on medical images. However, it is pointed out in (Yan et al., 2023) that the existing GPT4V is difficult to accurately analyze medical images and is prone to face the following major problems 1) Accurate localization requires clues. 2) Challenges in assessing object size. 3) Over-reliance on text. 4) Overemphasis on markers in images. 5) Very detailed answers. To mitigate errors in medical image analysis, we developed the Universal-Navigator Prompt (UNP) within IAI, which guides the MLLM in categorizing answers based on different types, such as

Exsample of Visual explanation generated from VQA-RAD



*Figure 9.* Example of Visual explanation generated with UNP from VQA-RAD.

questions and organs. In order to prevent GPT4V from being overly "confident" and giving incorrect inferences, we make full use of the information in the current conversation to label answers and give targeted answers based on different types of questions, such as questions and organs, which require more attention to confirming what type of modality is involved. For organ type we need to focus on the visual abnormality of the organ and we require GPT4 to reduce the focus on some image markers (*e.g.*, arrows) in the image, by which we greatly optimize GPT4V's ability to accurately analyze medical images. Since the dictionary information (Organ type, question type, etc.) of the SLAKE dataset is slightly different, we also made fine-tuning to UNP to adapt.

## A.2. Example of Visual explanation generated with UNP

As demonstrated in Figures 9 and 10, we show some of the results of generating visual interpretations on VQA-RAD and SLAKE dataset, using the IAI module. It can be found that for different types of organs (CHEST, HEAD, etc.) and different question types(PRES, POS, ORGAN etc.), IAI can make accurate and targeted visual interpretations. For example, for the POS type question "Where is the abnormality?", IAI responce the correct answer "left temporal lobe" .

## A.3. Qualitative Analysis of UNP's Impact on IAI

In addition, we further qualitatively compare the impact of UNP on MLLM in the IAI module, as shown in Figure 11. It can be found that GPT4V still lacks in understanding the details of medical images, location judgment and counting. By designing UNP, GPT4V has greatly improved its ability to accurately analyze medical images and can accurately generate

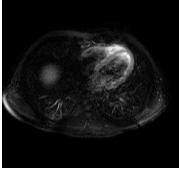Example of Visual explanation generated with UNP On SLAKE.



*Figure 10.* Example of Visual explanation generated with UNP On SLAKE.

visual explanations based on image and dialog information.

## A.4. UNP for Labeling Instruction Intent and Core Visual Explanation

Existing MLLMs have difficulty in understanding the "real intent" of instructions, focusing on irrelevant content, and capturing the visual information that users want to know from instructions. To this end, we use UNP to further mark the "real intent" of instructions and the visual evidence that can support the answer in the explanation. Then we construct a new IAI-Med VQA dataset with UNP, which can further help research in the field on Med-VQA. as shown in Figure 12. In order to guide MLLM in accurate labeling, UNP requires MLLM to focus more on medical terms, such as specific anatomical or pathological terms, and to be granular rather than labeling the entire problem, and to understand the "real intent" of the problem. Similarly, UNP has also made targeted designs for different organs and different types of problems. When marking the SLAKE dataset with "real_intention" and "core_visual_explanation", we also fine-tuned UNP to adapt to its data format.

## A.5. Example of Instruction Intent and Core Visual Explanation Labeling

The Figures 13 and 14, shows some examples of labeling on VQA-RAD and SLAKE datasets using UNP. After completing the generation of the visual explanations, in order to improve the model's instruction adherence and generalizability to unseen data, we used "real_intention" and "core_visual_explanation" to label the real intention of the question and the clues

*Figure 11.* Qualitative Analysis: With vs. Without UNP in IAI Module.

supporting the answer, respectively, and we can find that our UNP is able to recognize the core content of the question well, and identify the clues from the core content of the question, as well as the clues supporting the answer. It can be found that our UNP is able to recognize the core content of the question well and find the reason for the answer from the visual explanation, for example, our QA pairs with visual clues are enough to greatly improve the interpretability of MLLM on the Med-VQA task. We publicly label the IAI-Med VQA dataset to facilitate research in the field of Med VQA.

### A.6. Time Consumption for IAI-Med VQA Data Construction

In this part, we report the time required to construct the IAI-Med VQA dataset using the Universal-Navigator Prompt (UNP), measured in single GPU hours. The dataset construction process is divided into two main tasks: generating visual explanations and labeling instructions with "real intent" and "core visual explanations."

*Table 6.* GPU Time Consumption for Each Task in Dataset Construction

| Task | VQA-RAD (GPU hrs) | SLAKE (GPU hrs) |
|---|---|---|
| Generate visual explanation | 5 | 8 |
| Labeling instructions and explanation | 8 | 12 |

As shown in Table 6, for the VQA-RAD dataset, approximately 5 GPU hours were spent generating visual explanations and an additional 8 GPU hours for labeling tasks. In contrast, for the SLAKE dataset, the tasks required approximately 8 and 12 GPU hours, respectively.

Universal-Navigator prompt In IAI to label Instruction Intent and Core Visual Explanation

```
{
    "model": "gpt-4-vision-preview",
    "messages": [
        {
            "role": "system",
            "content": "You are an AI assistant specialized in biomedical topics. Given the latest question about a medical image and its
direct answer, your task is to use <real_intention></real_intention> tags to enclose the direct question to choose the most important
element to reflect the user's intentions. And <core_visual_explanation></core_visual_explanation> tags to wrap the description of the
visual content that supports the answer:
            1) The content should be related to the critical visual elements from the image and be as precise as possible,
            2) Each tag can be used repeatedly and the response should be complete, including the tags.
            3) Each tag can be used multiple times or once. Mark the question as many times as necessary. When marking answers, make sure
that the marked content includes complete visual descriptions.
            4) All question packages must be fine-grained using at least two tags. Pay attention to the marking of negative words!
            5) When labeling questions, label verbs and nouns rather than predicates. Don't mark combinations like 'are' and 'what is'
            6) When marking answers, try to mark the visual content that exists in the visual content of the image and supports the answer; and
the content is semantically complete and can be mapped to the corresponding individual in the image through text, and don't use more
than 2 for tagging.
             7) Identify and label the key medical terms in the question that are directly related to the diagnosis or condition being inquired
about, and focus your answer on these terms.
            8) Identify and label the specific anatomical or pathological terms present in the question.
            9) Identify and label the medical terminology within the question that relates to visible abnormalities or conditions in the image.
            10) It is not possible to wrap up all the content of the question, and important elements must be selected in a fine-grained manner:
Tagged content should have different focus points for different question types, such as:
-For MODALITY type questions, answers should focus on confirming what type of medical imaging the image is (e.g., X-ray, MRI, CT
scan, etc.).
 -For ORGAN type questions, answers should focus on the organs and their abnormalities visible in the image.
-For PRES type questions, the focus needs to be on confirming or ruling out the presence of specific pathological phenomena:
-For POS type questions, the focus is on describing the exact location where the anomaly was found.
For ABN type questions, the focus on explaining why a specific pathological abnormality is or is not present in the image..."
        },
        { "role": "user",
            "content": "The current question and answer of uploaded image is {question} and {answer}, The question type is {question
type} and the organ type is {image_organ}. You can refer to the reply methods in the template:{template}
            Here are some guidelines:
            1) When labeling questions, label verbs and nouns rather than predicates. Don't mark combinations like 'are' and 'what is':
            2) When marking answers, try to mark the visual content that exists in the visual content of the image and supports the answer,
and the content is semantically complete and can be mapped to the corresponding individual in the image through text, and don't use
more than 2 for tagging.
            3) To avoid wrapping invalid and redundant information, <real_intention></real_intention> must wrap a direct answer to the
question, which cannot exceed 10 words.
            4) <core_visual_explanation></core_visual_explanation> wrap the content cannot exceed 20 words. but please also return the
entire response! The contents of the package need to be as accurate as possible. "
        },{
            "type": "image_url",
            "image_url": {
                "url": "data:image/jpeg;base64,{base64_image}"
            }}],
    "max_tokens": 200}
```

*Figure 12.* UNP for Labeling Instruction Intent and Core Visual Explanation.

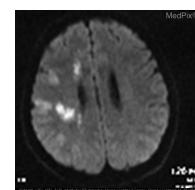Example of Instruction Intent and Core Visual Explanation Labeling On VQA-RAD

```
[
  {
    "conversations": [
      {
        "from": "human",
        "value": "Is there <real_intention>pericolic fat stranding</real_intention>?"
      },
      {
        "from": "gpt",
        "value": "Increased <core_visual_explain>attenuation around the colon, indicative of inflamed fat</core_visual_explain>, is consistent with pericolic fat stranding on this CT scan. The answer is yes"
      }
    ],
    "answer_type": "CLOSED",
    "image_organ": "ABD",
    "question_type": "PRES",
    "phrase_type": "freeform"
  },
  {
    "conversations": [
      {"from": "human", "value": "<real_intention>Where</real_intention> is the <real_intention>abnormality</real_intention> located?"
      },
      {
        "from": "gpt",
        "value": "Hyperdense area in the <core_visual_explain>right cerebellum</core_visual_explain>; indicative of a hemorrhage, lesion, or calcification. The answer is Right posteroinferior cerebellum. "
      }
    ],
    "answer_type": "OPEN",
    "image_organ": "HEAD",
    "question_type": "POS",
    "phrase_type": "frame"
  },]
```

```
[
  {
    "conversations": [
      {
        "from": "human",
        "value": "<real_intention>What organ system</real_intention> is <real_intention>pictured</real_intention>?"
      },
      {"from": "gpt", "value": "This is a brain MRI displaying <core_visual_explain>cerebral hemispheres and ventricular system</core_visual_explain> without skull structures. The answer is The brain. "
      }
    ],
    "answer_type": "OPEN",
    "image_organ": "HEAD",
    "question_type": "ORGAN",
    "phrase_type": "freeform"
  },
  "conversations": [
      {
        "from": "human",
        "value": "<real_intention>Which organ system</real_intention> is <real_intention>abnormal</real_intention> in this image?"
      },
      {
        "from": "gpt",
        "value": "The <core_visual_explain>cardiac silhouette is altered</task_related_visual>, a <core_visual_explain>cardiovascular system abnormality</core_visual_explain>. A <core_visual_explain>pacemaker device is visible</core_visual_explain>. The answer is Cardiovascular. "
      }
    ],
    "answer_type": "OPEN",
    "image_organ": "CHEST",
    "question_type": "ORGAN",
    "phrase_type": "freeform"
  },]
```

*Figure 13.* Example of Instruction Intent and Core Visual Explanation Labeling On VQA-RAD.

## A.7. Top 10 groupings by organ and question type on the VQA-RAD dataset

To ensure the interpretability of TC-Mix, we performed a grouping operation. Figure 15, shows the top10 sorting details grouped by organ type and problem type on the VQA-RAD dataset, which has 3 types of organs and more than 10 problem types. Since the visual cues corresponding to the same organ type and question type may be relatively similar and interrelated, this feature due to help the model to distinguish finer-grained differences. Therefore grouping Cut-Mix for different types can help MLLM to adapt to unseen features.

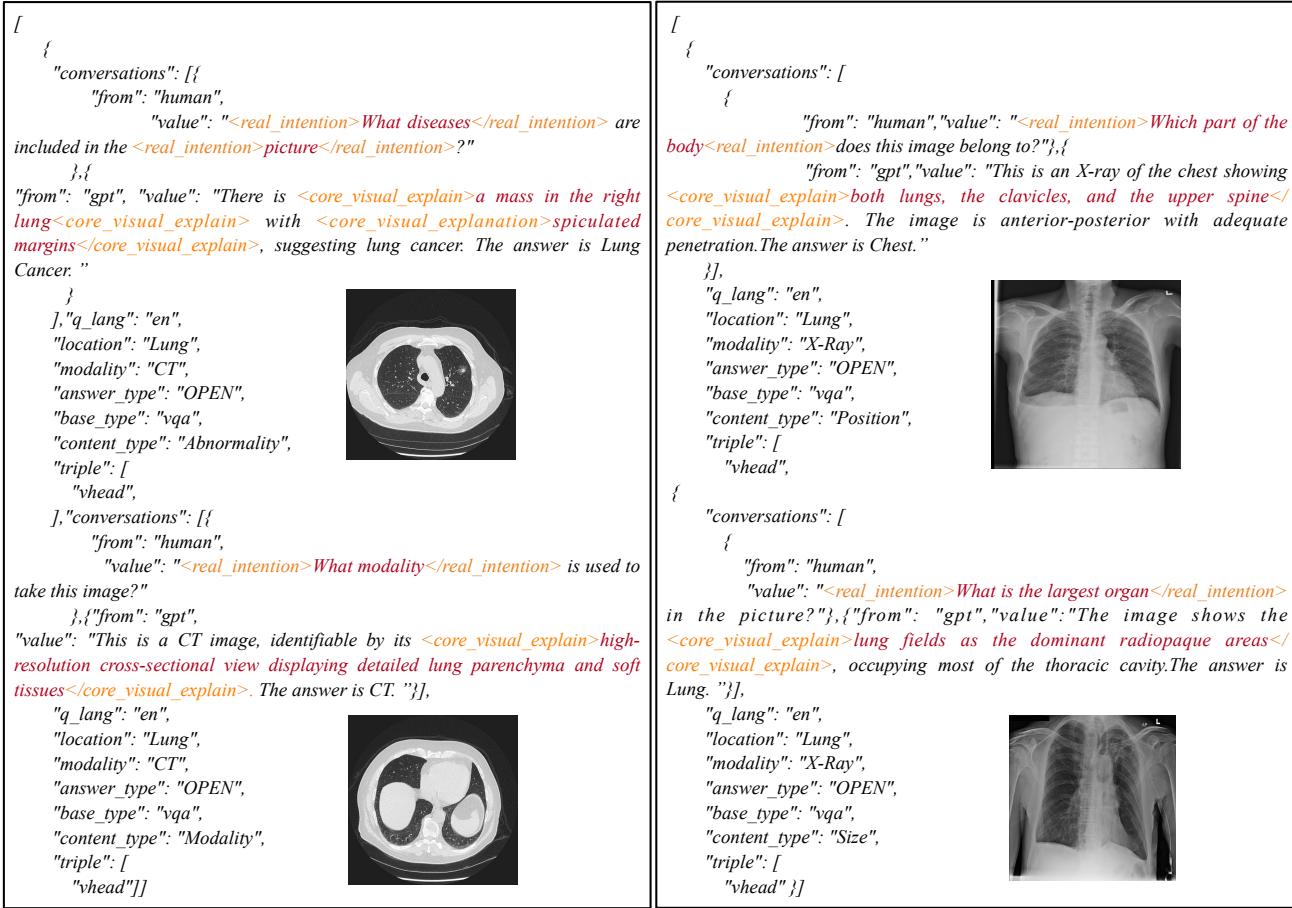Example of  Instruction Intent and Core Visual Explanation Labeling On SLAKE

```json
[
    {
      "conversations": [{
          "from": "human",
                    "value": "<real_intention>What diseases</real_intention> are
included in the <real_intention>picture</real_intention>?"
          },{
"from": "gpt", "value": "There is <core_visual_explain>a mass in the right
lung<core_visual_explain> with <core_visual_explanation>spiculated
margins</core_visual_explain>, suggesting lung cancer. The answer is Lung
Cancer. "
          }
      ],"q_lang": "en",
      "location": "Lung",
      "modality": "CT",
      "answer_type": "OPEN",
      "base_type": "vqa",
      "content_type": "Abnormality",
      "triple": [
        "vhead",
      ],"conversations": [{
          "from": "human",
             "value": "<real_intention>What modality</real_intention> is used to
take this image?"
          },{"from": "gpt",
"value": "This is a CT image, identifiable by its <core_visual_explain>high-
resolution cross-sectional view displaying detailed lung parenchyma and soft
tissues</core_visual_explain>. The answer is CT. "}],
      "q_lang": "en",
      "location": "Lung",
      "modality": "CT",
      "answer_type": "OPEN",
      "base_type": "vqa",
      "content_type": "Modality",
      "triple": [
        "vhead"]]
```

```json
[
  {
      "conversations": [
        {
              "from": "human","value": "<real_intention>Which part of the
body<real_intention>does this image belong to?"},{
              "from": "gpt","value": "This is an X-ray of the chest showing
<core_visual_explain>both lungs, the clavicles, and the upper spine</
core_visual_explain>. The image is anterior-posterior with adequate
penetration.The answer is Chest."
        }],
      "q_lang": "en",
      "location": "Lung",
      "modality": "X-Ray",
      "answer_type": "OPEN",
      "base_type": "vqa",
      "content_type": "Position",
      "triple": [
        "vhead",
  {
      "conversations": [
        {
          "from": "human",
          "value": "<real_intention>What is the largest organ</real_intention>
in the picture?"},{"from": "gpt","value":"The image shows the
<core_visual_explain>lung fields as the dominant radiopaque areas</
core_visual_explain>, occupying most of the thoracic cavity.The answer is
Lung. "}],
      "q_lang": "en",
      "location": "Lung",
      "modality": "X-Ray",
      "answer_type": "OPEN",
      "base_type": "vqa",
      "content_type": "Size",
      "triple": [
        "vhead" }]
```

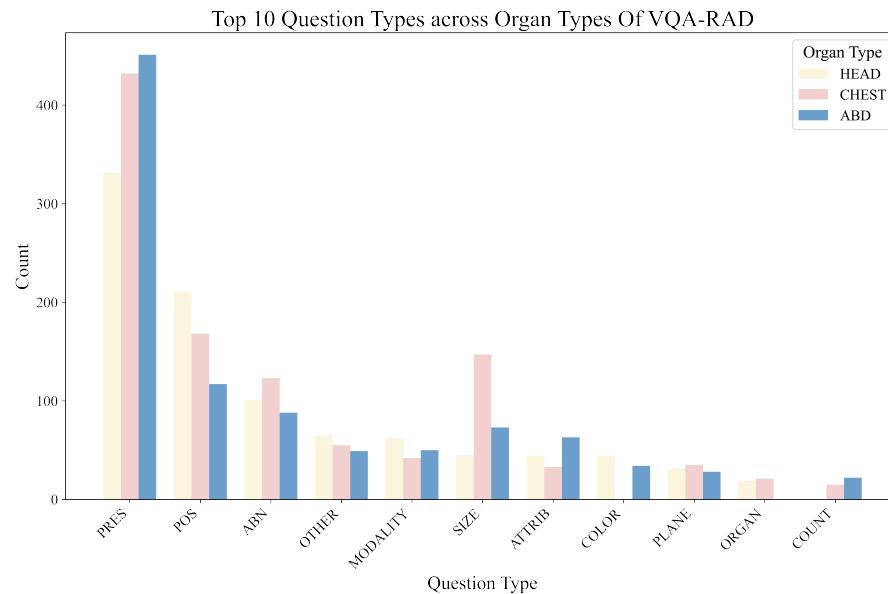*Figure 14.* Example of Instruction Intent and Core Visual Explanation Labeling On SLAKE.



*Figure 15.* Top 10 groupings by organ and question type on the VQA-RAD dataset.