

## 7. Supplementary

### 7.1. Detail related work

**Learnable number embedding.** A straightforward approach for incorporating numerical values is using learnable networks to continuously project numerical values into and out of the  $d$ -dimensional embedding space. This approach typically involves two small neural networks: a number encoder and a number decoder.

Specifically, the encoding function  $\mathcal{G}_{enc}$  uses a number encoder,  $f_{enc} : \mathbb{R} \rightarrow \mathbb{R}^d$ , which maps the input number  $m$  to a  $d$ -dimensional vector. This vector is then combined with the base embedding for the [NUM] token, for instance, through addition. The decoding function  $\mathcal{G}_{dec}$  uses a number decoder,  $f_{dec} : \mathbb{R}^d \rightarrow \mathbb{R}$ , which performs the reverse operation: it takes the final hidden state from the model and projects it back to a scalar numerical value. Both  $f_{enc}$  and  $f_{dec}$  are typically implemented as Multi-Layer Perceptrons (MLPs). The mapping functions are thus defined as:

$$\begin{aligned} \mathcal{G}_{enc}(\mathbf{x}_{[\text{NUM}]}, m) &= \mathbf{x}_{[\text{NUM}]} + f_{enc}(m) \\ \mathcal{G}_{dec}(\hat{\mathbf{x}}_i) &= f_{dec}(\hat{\mathbf{x}}_i), \\ &\text{if } [\text{NUM}] = \arg \min_v f_{lm}(\hat{\mathbf{x}}_i), \end{aligned} \quad (15)$$

where  $f_{lm}$  refers to original language modeling head from LLM.

Although this approach is notable for its simplicity as it parametrizes both mapping functions with learnable MLPs and can theoretically satisfies **C1** requirement of "wide-range of value" and **C3**, it introduces two critical challenges. First, an MLP does not inherently guarantee that the ordinal relationship of numbers is preserved in the high-dimensional embedding space. There is no inductive bias to ensure that the embedding for a number (e.g. the number 3 is closer to 4 than to 30. I.e., the embedding  $f_{enc}(3)$  is not guaranteed to be closer to  $f_{enc}(4)$  than it is to  $f_{enc}(30)$ ). The learned geometric relationships can be arbitrary without additional constraints, which makes **C1** "ordinal-preserving" requirement difficult to satisfy. Second, and more fundamentally, the projection network suffers from poor out-of-distribution (OOD) generalization and prone to context noise. An MLP trained on numbers within a specific range cannot be expected to produce meaningful embeddings for numbers it has not seen during training. Furthermore, during decoding, the MLP may struggle to accurately reconstruct the original number from the hidden state, especially if the hidden state has been influenced by other contextual information. Therefore, this approach cannot reliably satisfy part of **C1** and **C2**, making it brittle for applications with unbounded numerical data.

**Magnitude-based scaling.** The xVal [13] framework injects numerical values by scaling the embedding magnitude of the special token. To encode a numerical value  $m_i$ ,

the method multiplies the base embedding  $\mathbf{x}_{[\text{NUM}]}$  with a normalized version of the value, such that the norm of the resulting embedding scales with the magnitude of  $m_i$ . During training, a separate regression head  $f_{num}$  is trained with a Mean-Squared-Error (MSE) loss to predict the original value  $m_i$  from the final hidden state. During inference, a standard language model head  $f_{lm}$  predicts the next token. If it predicts [NUM], the number head  $f_{num}$  is triggered to generate the numerical value  $\hat{m}_i$  from the corresponding hidden state  $\hat{\mathbf{x}}_i$ . The  $\mathcal{G}_{enc}$  and  $\mathcal{G}_{dec}$  functions for this approach are defined as:

$$\begin{aligned} \mathcal{G}_{enc}(\mathbf{x}_{[\text{NUM}]}, m_i) &= \tanh(m_i) \cdot \mathbf{x}_{[\text{NUM}]} \\ \mathcal{G}_{dec}(\hat{\mathbf{x}}_i) &= f_{num}(\hat{\mathbf{x}}_i), \\ &\text{if } [\text{NUM}] = \arg \min_v f_{lm}(\hat{\mathbf{x}}_i) \end{aligned} \quad (16)$$

Although xVal guarantees an **C1** ordinal relationship by directly scaling the embedding with the number's magnitude, it suffers from critical limitations. First, it cannot encode an unbounded range of numbers, as it requires a normalization function like  $\tanh$  on the input number to prevent the embedding norm from exploding. This leads to value saturation, where the model cannot distinguish between different large numbers and fail for **C1** "wide-range value" requirement. More importantly, xVal was originally tested on Pre-LayerNorm (LN) [2] architectures like GPT-2 [34]. However, all recent LLMs turns to Post-LN due to training difficulties [42]. As pointed out by Brody et al. [6], the norm of a vector after going through LN is always  $\sqrt{d}$ , which diminished xVal's scaling utility to modern LLM architecture and failed at **C3**.

### 7.2. Algorithm outline

Here we describe the algorithm outlines for Rotary Number Encoding (RNE) in Alg. 1 and Decoding (RND) in Alg. 2.

---

#### Algorithm 1 Rotary Number Encoding (RNE)

---

**Input:**  $\mathbf{x}_{[\text{NUM}]} \in \mathbb{R}^d$ , value  $m$ .

**Output:** Encoded  $\mathbf{x}^m \in \mathbb{R}^d$ .

- 1: **procedure** RNE( $\mathbf{x}_{[\text{NUM}]}, m$ )
  - 2:    $\Theta : \theta_j \leftarrow B^{-2j/d}$ , for  $j \in \{0, \dots, d/2 - 1\}$ .
  - 3:   2D Rotation:  $\mathbf{R}_2(m\theta_j) \leftarrow \begin{bmatrix} \cos m\theta_j & -\sin m\theta_j \\ \sin m\theta_j & \cos m\theta_j \end{bmatrix}$ .
  - 4:    $d$ -dim Rotation:  $\mathbf{R}_d(m\Theta) \leftarrow \text{diag}(\mathbf{R}_2(m\theta_j))$ .
  - 5:   **return**  $\mathbf{R}_d(m\Theta)\mathbf{x}_{[\text{NUM}]}$ .
- 

### 7.3. Distinguishability–smoothness trade-off

We analyze how the weighting exponent  $p$  affects the geometry of the scalar score  $\mathcal{S}(m, p)$ . Let  $B > 1$ , even  $d$ , and

$$\begin{aligned} \theta_j &:= B^{-2j/d}, \quad j = 1, \dots, d/2, \\ w_j(p) &:= \theta_j^{-p} \|\mathbf{x}_{[\text{NUM}]_j}\|^2. \end{aligned} \quad (17)$$

---

**Algorithm 2** Rotary Number Decoding (RND)
 

---

**Input:** Model output  $\hat{\mathbf{x}}_i \in \mathbb{R}^d$ , candidates  $\mathcal{M}$ , parameter  $p$ .

**Output:** Decoded number  $\hat{m}$ .

- 1: **procedure** RND( $\hat{\mathbf{x}}_i, \mathcal{M}, p$ )
  - 2:    $\Theta : \theta_j \leftarrow B^{-2j/d}$  for  $j \in \{1, \dots, d/2 - 1\}$ .
  - 3:   Compute input signal:  $\hat{s} \leftarrow \mathcal{S}(\hat{\mathbf{x}}_i, p)$ .
  - 4:   Compute candidate targets:  $s'_{m'} \leftarrow \mathcal{S}(m', p)$  for each  $m' \in \mathcal{M}$ .
  - 5:   Predict:  $\hat{m} \leftarrow \arg \min_{m' \in \mathcal{M}} (\hat{s} - s'_{m'})^2$ .
  - 6:   **return**  $\hat{m}$ .
- 

Then

$$\mathcal{S}(m, p) := \sum_{j=1}^{d/2} w_j(p) \cos(m\theta_j). \quad (18)$$

**Lemma 7.1** (Derivative and oscillation scale). *The derivative of the score with respect to  $m$  is*

$$\begin{aligned} \mathcal{S}'(m, p) &= - \sum_{j=1}^{d/2} w_j(p) \theta_j \sin(m\theta_j) \\ &= - \sum_{j=1}^{d/2} \|\mathbf{x}_{[\text{NUM}],j}\|^2 \theta_j^{1-p} \sin(m\theta_j). \end{aligned} \quad (19)$$

Each term oscillates with period  $T_j = 2\pi/\theta_j$ . For  $p < 1$ , the factor  $\theta_j^{1-p}$  increases with  $\theta_j$ , so high-frequency components ( $\theta_j$  large, short  $T_j$ ) dominate  $\mathcal{S}'$ ; for  $p > 1$ ,  $1 - p < 0$  and low-frequency components ( $\theta_j$  small, long  $T_j$ ) dominate.

Lemma 7.1 shows that increasing  $p$  shifts  $\mathcal{S}'(m, p)$  from high- to low-frequency components, making  $\mathcal{S}(m, p)$  smoother in  $m$ .

**Phase assumption.** We now quantify typical changes in the score between nearby values of  $m$ . We assume that, as  $m$  ranges over a large interval, the phases  $m\theta_j \bmod 2\pi$  explore  $[0, 2\pi)$  approximately uniformly, so that

$$\mathbb{E}_m[\sin(m\theta_j)] \approx 0, \quad \mathbb{E}_m[\cos(m\theta_j)] \approx 0, \quad \mathbb{E}_m[\sin^2(m\theta_j)] \approx \frac{1}{2}. \quad (20)$$

**Lemma 7.2** (Typical dot-product gap and RMS scale). *Let  $\Delta\mathcal{S}_\delta(m, p) := \mathcal{S}(m + \delta, p) - \mathcal{S}(m, p)$ . Under the phase assumption,*

$$\mathbb{E}_m[\Delta\mathcal{S}_\delta(m, p)^2] \approx 2 \sum_{j=1}^{d/2} w_j(p)^2 \sin^2\left(\frac{\theta_j \delta}{2}\right), \quad (21)$$

and for  $\theta_j |\delta| \ll 1$  for all  $j$ ,

$$\begin{aligned} \Delta_{\text{typ}}(\delta, p) &:= \sqrt{\mathbb{E}_m[\Delta\mathcal{S}_\delta(m, p)^2]} \\ &\approx |\delta| C(p), \quad C(p) := \left(\frac{1}{2} \sum_{j=1}^{d/2} w_j(p)^2 \theta_j^2\right)^{1/2}. \end{aligned} \quad (22)$$

Moreover, the typical scale of the score is

$$\sigma_S(p)^2 := \mathbb{E}_m[\mathcal{S}(m, p)^2] \approx \frac{1}{2} \sum_{j=1}^{d/2} w_j(p)^2. \quad (23)$$

We measure distinguishability between consecutive integers by the *relative dot-product gap*:

**Definition 1** (Relative dot-product gap). *For consecutive integers ( $\delta = 1$ ), define*

$$\begin{aligned} r(p) &:= \frac{\Delta_{\text{typ}}(1, p)}{\sigma_S(p)} \\ &\approx \sqrt{\frac{\sum_{j=1}^{d/2} w_j(p)^2 \theta_j^2}{\sum_{j=1}^{d/2} w_j(p)^2}} \\ &= \sqrt{\frac{\sum_{j=1}^{d/2} \|\mathbf{x}_{[\text{NUM}],j}\|^4 \theta_j^{2(1-p)}}{\sum_{j=1}^{d/2} \|\mathbf{x}_{[\text{NUM}],j}\|^4 \theta_j^{-2p}}}. \end{aligned} \quad (24)$$

**Proposition 1** (Distinguishability–smoothness trade-off).

Assume at least two distinct  $\theta_j$  have  $\|\mathbf{x}_{[\text{NUM}],j}\| > 0$ . Then:

1. Increasing  $p$  shifts  $\mathcal{S}'(m, p)$  towards lower frequencies (Lemma 7.1), so  $\mathcal{S}(m, p)$  becomes smoother in  $m$ .
2. The relative gap  $r(p)$  is strictly decreasing in  $p$ . Hence larger  $p$  produces smoother but less distinguishable scalar scores, and smaller  $p$  produces more oscillatory but more separated scores between neighboring values.

*Proof.* (1) follows directly from Lemma 7.1. For (2), write

$$\begin{aligned} r(p)^2 &= \frac{N(p)}{D(p)}, \\ N(p) &= \sum_{j=1}^{d/2} a_j \theta_j^{2(1-p)}, \\ D(p) &= \sum_{j=1}^{d/2} a_j \theta_j^{-2p}, \\ a_j &:= \|\mathbf{x}_{[\text{NUM}],j}\|^4 > 0. \end{aligned} \quad (25)$$

Differentiating and symmetrizing over  $(j, k)$  gives

$$\begin{aligned} \frac{d}{dp} r(p)^2 &= \frac{N'(p)D(p) - N(p)D'(p)}{D(p)^2} \\ &= - \sum_{j < k} a_j a_k (\ln \theta_j - \ln \theta_k) (\theta_j^2 - \theta_k^2) \theta_j^{-2p} \theta_k^{-2p}. \end{aligned} \quad (26)$$

Since the  $\theta_j$  are strictly decreasing, each term is negative whenever  $a_j, a_k > 0$ , so  $r(p)^2$  and hence  $r(p)$  are strictly decreasing in  $p$ .  $\square$

#### 7.4. Precision-dependent choice of weighting $p$

We now consider a practical constraint on the choice of  $p$  imposed by the numerical precision of the floating-point format used to represent  $\mathcal{S}(m, p)$ . In modern LLMs, full precision (fp32) is often not used for inference due to computational and memory constraints. Instead, lower-precision formats such as fp16, bfloat16, or even fp8 are commonly employed. These formats have limited precision, which can lead to quantization errors when representing the scalar scores  $\mathcal{S}(m, p)$ .

We start relate the distinguishability measure  $r(p)$  to the numerical precision of the floating-point format used to represent  $\mathcal{S}(m, p)$ . Recall that  $\Delta_{\text{typ}}(1, p)$  denotes the typical change in the score between neighboring integers  $m$  and  $m + 1$ , and  $\sigma_S(p)$  denotes the typical scale (RMS) of the score. Robust scalar decoding requires that the score change between neighbors is large compared to the quantization step of the format.

**Unit roundoff and quantization step.** To quantify the limits of numerical precision, we consider the **unit roundoff**, denoted by  $u$ . For a standard base-2 floating-point format with  $n_{\text{sig}}$  significand bits (including the implicit leading bit), the resolution is determined by the distance between 1.0 and next representable floating-point number. This distance is given by:

$$u = 2^{-(n_{\text{sig}}-1)}. \quad (27)$$

This value  $u$  represents the relative quantization error. Consequently, for a scalar score with a typical magnitude  $\sigma_S(p)$ , the absolute spacing between representable values is approximately  $u \cdot \sigma_S(p)$ .

In our setting, we want the typical separation between the scores  $m$  and  $m + 1$  to be strictly resolvable, that is much larger than this quantization step:

$$\Delta_{\text{typ}}(1, p) \gg u, \sigma_S(p) \leftrightarrow r(p) := \frac{\Delta_{\text{typ}}(1, p)}{\sigma_S(p)} \gg u. \quad (28)$$

**Theorem 7.3** (Precision-dependent constraint on  $p$ ). *For each floating-point format (fixed  $u$ ), there exists a unique  $p_{\text{crit}}(u)$  such that*

$$r(p_{\text{crit}}(u)) = u, \quad (29)$$

with

$$\begin{aligned} p < p_{\text{crit}}(u) &\Rightarrow r(p) > u \quad (\text{neighbors well-separated}), \\ p > p_{\text{crit}}(u) &\Rightarrow r(p) < u \quad (\text{neighbors can collapse}). \end{aligned} \quad (30)$$

Combined with Proposition 1, this yields a three-way trade-off between smoothness, distinguishability, and numerical precision.

*Proof.* By Proposition 1,  $r(p)$  is continuous and strictly decreasing in  $p$ . Hence, for any  $u$  in the range of  $r$ , there is a unique solution of  $r(p) = u$ , which we denote  $p_{\text{crit}}(u)$ . The implications for  $p < p_{\text{crit}}(u)$  and  $p > p_{\text{crit}}(u)$  follow directly from the inequalities  $r(p) \gtrless u$ .  $\square$

**Corollary 1** (Recommended  $p$  for common precisions (example)). *For common IEEE-like formats, the effective spacing  $u$  is determined by the number of significand bits  $n_{\text{sig}}$ . The relationship is given by  $u \approx 2^{-(n_{\text{sig}}-1)}$ . For typical formats:*

- *fp32:*  $n_{\text{sig}} = 24 \Rightarrow u \approx 2^{-23} \approx 1.2 \times 10^{-7}$ ,
- *fp16:*  $n_{\text{sig}} = 11 \Rightarrow u \approx 2^{-10} \approx 10^{-3}$ ,
- *bfloat16:*  $n_{\text{sig}} = 8 \Rightarrow u \approx 2^{-7} \approx 8 \times 10^{-3}$ ,
- *fp8 variants:*  $n_{\text{sig}} \approx 3-4 \Rightarrow u \approx O(10^{-2} - 10^{-1})$ .

*For a typical configuration ( $B = 5 \cdot 10^5$ ,  $d = 4096$ ) and roughly uniform block energies, numerical evaluation of  $r(p)$  with these  $u$  values yields:*

$$\begin{aligned} \text{fp32:} \quad & r(p) \gg u \quad \forall p \text{ (no practical constraint);} \\ \text{fp16:} \quad & r(p) \gg u \text{ for } p \lesssim 0.3; \\ \text{bfloat16:} \quad & r(p) \gg u \text{ for } p \lesssim 0.2; \\ \text{fp8:} \quad & r(p) \lesssim u \quad \forall p \text{ (scalar score insufficient)}. \end{aligned} \quad (31)$$

**Conclusion:** Thus, for scalar dot-product decoding, small  $p$  (e.g.,  $p \in [0, 0.3]$ ) is critical when using lower-precision formats like fp16 and bfloat16 to ensure the score separation exceeds the quantization noise. Conversely, fp32 imposes essentially no restriction, while fp8 possesses insufficient resolution for scalar decoding regardless of  $p$ , necessitating either full-vector decoding or higher-precision accumulation. This theoretical analysis justifies our ablation results in Tab. 4 and explains why  $p \approx 0.2 - 0.3$  is optimal for optimal operating point for standard inference hardware.

#### 7.5. Dataset curation

We curated a dataset with 25,645 non-contrast Chest computed tomography (CT) images from 14,218 unique patients, each with its corresponding radiology reports. The radiology report contains a "Findings" section and an "Impression" section. The "Findings" section describes the detailed observations from the CT scan, while the "Impression" section provides a concise summary of the key findings and their clinical significance. For dataset splitting, we follow a 80-20 patient-wise approach to ensure that the training and testing sets are completely disjoint. We finally have 19,944 training samples with 11,059 patients and 5,701 testing samples with 3,159 patients.

To create a comprehensive training dataset for our model, we employed a two-stage generation process using two different LLMs, followed by CT-RATE [15] and HuatuoGPT [9] curation process. First, we used Llama-3.1 [14] to generate basic VQA pairs (multiple-choice, short-answer, summarization) for information extraction. Second, to enhance diversity, we used Gemini-2.5 [38] to simulate various clinical scenarios to generating questions with different tones and formats. The comprehensive generation process covering both simple queries and detailed, long-form responses, improving model’s instruction following capabilities. The full generation process and examples are illustrated in Fig. 6.

For evaluation on radiology measurement estimation task, we combined traditional regular expression filtering and LLM based filtering to locate label of interested and their measurements. We first use regular expression to extract numerical values and their related text from the raw radiology report. Then we use LLM to filter out low-frequency labels and identify them as label of interested. The selected labels are: **Pulmonary Nodule, Ascending aorta, Aortic root, Aortic arch, Descending aorta, Aortic hiatus, Pulmonary Artery, Lymph Node, Thyroid nodule, Renal cyst and Adrenal lesion**. Then we use LLM to extract additional context related to the numerical value, forming meta-data query consists of number, unit, label and context. Finally, we use LLM to generate QA pairs based on meta-data. The full pipeline is illustrated in Fig. 5.

## 7.6. Medical Grounding Evaluation Details

We presented the full table from MedSeq-Bench [45] at Tab. 7. The table includes the performance of various general-purpose MLLMs across different grounding tasks, measured by Interaction-Over-Union (IOU) and Acc@5 metrics. IOU quantifies the overlap between the predicted bounding box and the ground truth bounding box. Acc@5 measure the proportion of predictions whose IOU with the ground truth exceeds 0.5. The results are categorized into two main sections: Image Difference Grounding and Image Consistency Grounding, with specific tasks listed under each category. The average performance across all tasks is also provided for each model.

## 7.7. Limitations

The score-matching decoder searches over a fixed candidate set  $\mathcal{M}$ , reintroducing a form of discretization; however, the step size of  $\mathcal{M}$  is configurable per task and training uses a continuous regression loss. For large coordinate ranges (e.g., 0–1024 in high-resolution grounding), aliasing may occur at the highest rotational frequencies; this can be mitigated by adjusting  $p$  or adopting frequency scaling strategies analogous to YaRN. Currently, RNED supports only non-negative values; handling signed quantities (e.g., Hounsfield units) would require a sign-magnitude decomposition, which we

leave as future work. Our primary evaluation uses the in-house Opport-CT dataset; broader validation on additional public benchmarks is an important future direction. Finally, RNE is applied only at the input embedding layer, without an explicit guarantee that intermediate transformer layers preserve the rotational structure; nonetheless, the end-to-end regression loss ensures that the output embedding carries sufficient numerical information for accurate decoding.

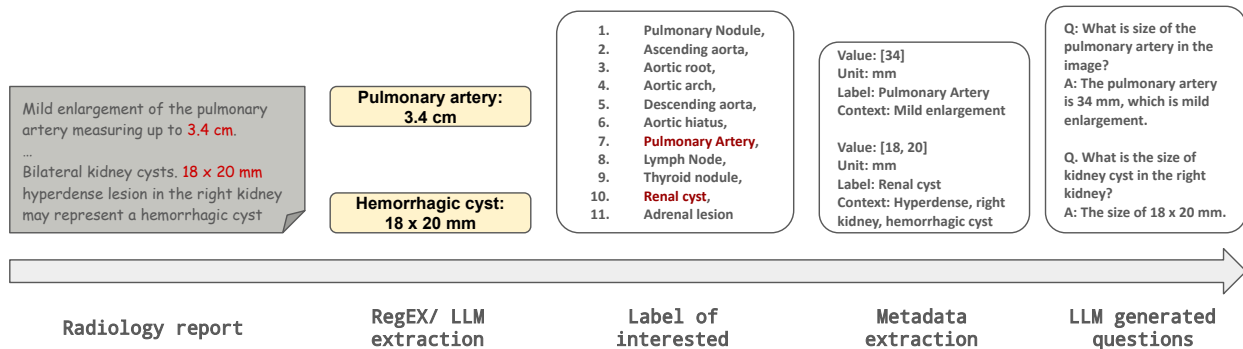


Figure 5. Pipeline for generating numerical decoding testing set. 1. We first using regular expression and LLM to extract numerical value and its related text from raw Radiology report. 2. Based on the extracted text, we conclude the 11 interested labels and remove low-frequency texts. 3. We extract additional context related to numerical value, forming meta-data query consists of number, unit, label and context. 4. We use LLM to generate QA pairs based on meta-data.

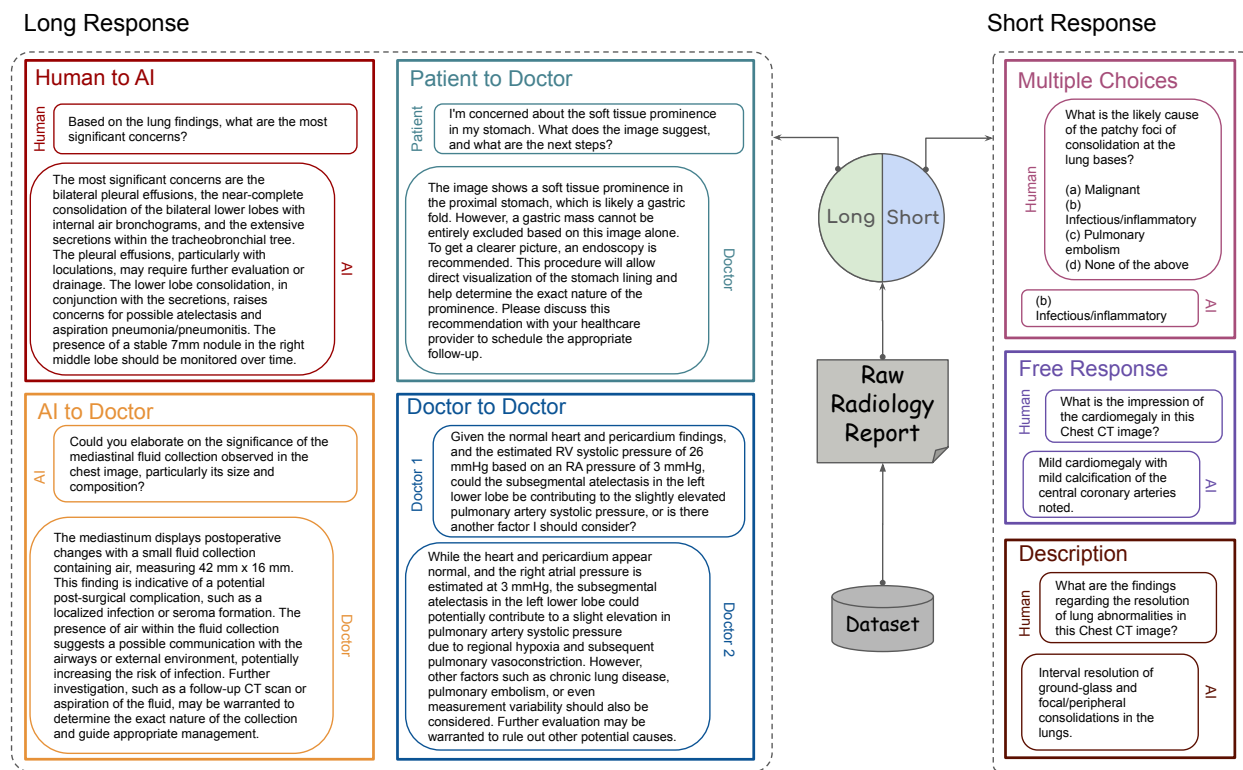


Figure 6. Sample figure illustrating the training data generation. Each Radiology report is augmented to generate multiple VQA examples that represent different perspectives. On the left, we have several examples of “long responses” that use a detailed description under different prompt scenarios that make model generated results be detailed and diverse. On the right, we show shorter question examples with multiple choice answers, a brief description and a free response that focuses on a specific information directly extracted from the report.

Model	Size	Metric	Image Difference Grounding		Image Consistency Grounding					Avg.	
			Reg. DG	Non-reg. DG	Multi-view	Obj. Track.	Concept	Patch	Cross-modal		Referring
<b>General-purpose MLLMs</b>											
Qwen2.5-VL [3]	3B	IoU Acc@5	0.59 0.30	1.62 1.30	7.12 3.90	21.32 16.80	6.98 0.80	27.36 3.40	10.02 1.65	12.99 6.82	10.94 4.20
Qwen2.5-VL [3]	7B	IoU Acc@5	0.88 0.30	1.25 0.00	8.48 3.73	22.41 17.80	4.22 1.00	28.87 5.70	16.29 4.45	12.58 6.21	12.31 4.90
Qwen2.5-VL [3]	32B	IoU Acc@5	2.69 1.40	3.48 1.20	7.35 2.61	19.12 13.40	6.53 1.30	26.92 7.10	12.59 4.90	18.71 11.67	12.47 5.71
Qwen2.5-VL [3]	72B	IoU Acc@5	4.37 2.60	3.46 0.80	7.22 2.78	13.11 7.70	10.33 3.50	26.45 6.30	16.32 7.00	20.19 14.10	13.35 6.12
MiniCPM-V-2.6 [43]	8B	IoU Acc@5	1.36 0.00	1.50 0.00	15.82 5.20	24.03 18.50	9.90 2.10	28.65 12.20	12.72 3.30	12.44 3.64	13.24 5.27
MiniCPM-O-2.6 [43]	8B	IoU Acc@5	1.69 0.10	1.63 0.00	12.11 2.43	15.25 9.60	9.88 1.70	22.96 9.20	9.53 2.35	8.82 2.02	10.12 3.23
mPLUG-Owl3 [44]	7B	IoU Acc@5	2.12 0.00	2.55 0.00	15.64 3.64	15.62 4.40	6.80 0.80	30.42 3.60	17.06 4.80	11.92 5.47	13.22 3.19
Mantis-Idefics2 [20]	8B	IoU Acc@5	0.49 0.00	0.62 0.00	18.69 8.59	28.04 23.50	6.27 0.50	10.26 1.10	9.59 0.95	6.05 0.54	9.90 3.91
LLaVA-OneVision [21]	7B	IoU Acc@5	1.09 0.00	0.01 0.00	9.26 1.13	10.50 3.20	11.33 1.80	22.20 5.30	19.08 6.70	17.11 5.67	12.39 3.47
LLaVA-OneVision [21]	72B	IoU Acc@5	2.58 0.80	2.87 0.90	11.74 1.39	9.61 2.30	10.95 3.30	32.38 20.30	16.24 5.40	15.43 6.68	13.21 5.18
InternVL2 [11]	8B	IoU Acc@5	0.18 0.00	0.38 0.00	17.34 7.03	26.45 21.20	5.56 0.80	10.36 0.70	6.23 1.00	15.73 7.69	10.24 4.59
InternVL2 [11]	76B	IoU Acc@5	0.15 0.00	0.15 0.00	10.00 3.90	15.56 11.80	3.39 0.40	6.64 1.10	2.83 0.75	15.69 9.92	6.88 3.53
InternVL2.5 [10]	8B	IoU Acc@5	0.26 0.00	0.38 0.00	13.52 3.56	20.82 13.80	1.96 0.00	5.25 0.00	4.70 0.85	9.56 3.44	7.04 2.56
InternVL2.5 [10]	78B	IoU Acc@5	0.24 0.10	0.32 0.10	9.16 2.08	16.18 10.00	4.32 0.50	11.86 2.30	5.48 1.25	10.67 4.52	7.29 2.55
InternVL3 [48]	8B	IoU Acc@5	1.07 0.30	1.20 0.00	14.36 4.42	13.30 6.50	6.43 0.90	18.73 4.60	4.73 1.15	15.16 7.42	9.26 3.19
InternVL3 [48]	14B	IoU Acc@5	0.66 0.00	0.71 0.00	13.24 5.31	19.77 13.00	8.60 2.10	13.17 2.40	10.87 3.70	14.57 7.76	10.53 4.41
InternVL3 [48]	38B	IoU Acc@5	0.98 0.10	1.76 0.20	12.99 4.79	19.27 13.60	7.63 2.10	17.76 2.90	6.47 1.75	16.59 10.05	10.37 4.44
InternVL3 [48]	78B	IoU Acc@5	0.20 0.00	0.53 0.00	6.35 2.43	13.03 8.00	3.57 0.90	11.81 2.50	3.34 0.85	12.76 8.10	6.44 2.90
Migician [24]	7B	IoU Acc@5	15.26 7.80	14.49 6.10	18.16 7.84	21.38 14.90	14.23 7.20	28.87 13.70	21.41 12.15	25.30 18.02	20.29 11.39
<b>Medical-domain specialized MLLMs</b>											
HuatuoGPT-Vision [9]	7B	IoU Acc@5	1.35 0.00	1.84 0.20	10.42 2.78	14.57 9.20	7.99 0.80	15.52 2.30	9.46 2.15	9.60 1.82	8.97 2.36
HuatuoGPT-Vision [9]	34B	IoU Acc@5	1.44 0.00	2.15 0.00	9.41 1.65	13.25 8.30	6.43 0.70	14.53 1.40	10.60 2.60	8.60 1.75	8.57 2.09
MedSeq-Grounder [45]	7B	IoU Acc@5	<b>83.29</b> <b>93.20</b>	<b>83.72</b> <b>94.10</b>	55.03 60.19	62.10 67.20	74.11 82.60	<b>85.25</b> <b>98.80</b>	78.77 82.75	60.43 65.59	72.55 79.71
MedSeq-Grounder + RNED	7B	IoU Acc@5	82.83 92.30	81.13 91.95	<b>59.26</b> <b>65.06</b>	<b>63.15</b> <b>68.60</b>	<b>75.37</b> <b>83.52</b>	84.11 96.80	<b>83.54</b> <b>87.00</b>	<b>64.29</b> <b>71.70</b>	<b>74.21</b> <b>82.33</b>

Table 7. Performance on image difference and image consistency grounding tasks on MedSG-Bench. We report IoU and Acc@5.