

# Instruction-Guided Lesion Segmentation for Chest X-rays with Automatically Generated Large-Scale Dataset

Anonymous CVPR submission

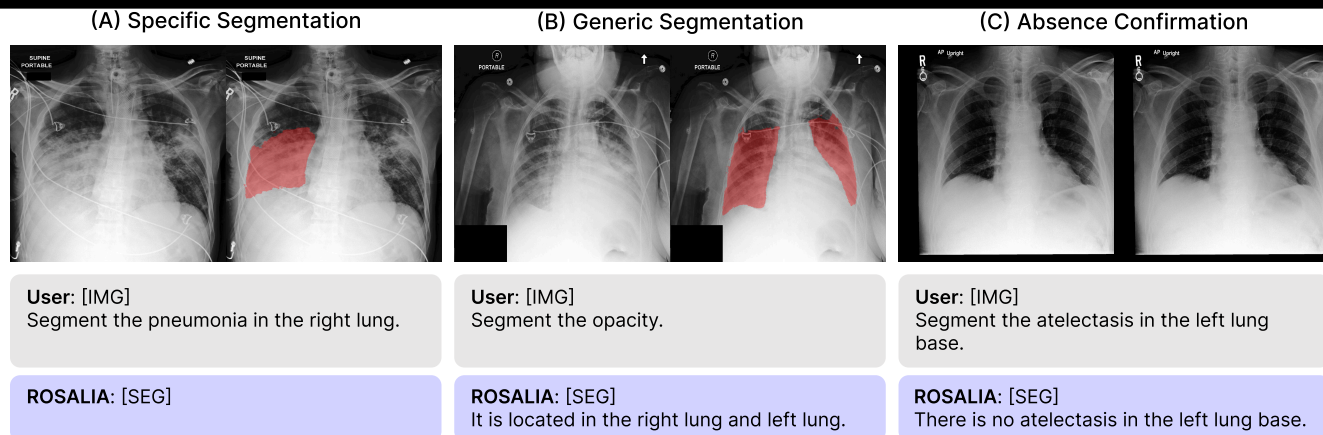


Figure 1. Examples of the instruction-guided CXR lesion segmentation task. Given text instructions for various lesion types and locations of interest, ROSALIA, a VLM trained on our MIMIC-ILS dataset, can: (A) segment lesions in a specified location, (B) segment lesions globally, and (C) detect empty-target cases. As can be seen in (A), ROSALIA correctly ignores the unrequested lesion in the left lung.

## Abstract

001 *The applicability of current lesion segmentation models for*  
002 *chest X-rays (CXRs) has been limited both by a small number*  
003 *of target labels and the reliance on long, detailed expert-level*  
004 *text inputs, creating a barrier to practical use. To address*  
005 *these limitations, we introduce a new paradigm: instruction-*  
006 *guided lesion segmentation (ILS), which is designed to seg-*  
007 *ment diverse lesion types based on simple, user-friendly*  
008 *instructions. Under this paradigm, we construct MIMIC-ILS,*  
009 *the first large-scale instruction-answer dataset for CXR le-*  
010 *esion segmentation, using our fully automated multimodal*  
011 *pipeline that generates annotations from chest x-ray im-*  
012 *ages and their corresponding reports. MIMIC-ILS contains*  
013 *1.1M instruction-answer pairs derived from 192K images*  
014 *and 91K unique segmentation masks, covering seven major*  
015 *lesion types. To empirically demonstrate its utility, we in-*  
016 *troduce ROSALIA, a vision-language model fine-tuned on*  
017 *MIMIC-ILS. ROSALIA can segment diverse lesions and pro-*  
018 *vide textual explanations in response to user instructions.*  
019 *The model achieves high segmentation and textual accuracy*  
020 *in our newly proposed task, highlighting the effectiveness of*  
021 *our pipeline and the value of MIMIC-ILS as a foundational*  
022 *resource for pixel-level CXR lesion grounding.*

## 1. Introduction

Medical imaging is an essential technique in modern  
medicine, enabling accurate diagnosis and appropriate treat-  
ment. Among various imaging modalities, chest X-ray (CXR)  
is one of the most common examinations due to its high ac-  
cessibility and rapid acquisition [4]. Radiologists reach a  
diagnosis by integrating visual evidence from CXRs with  
their clinical knowledge, and describe these findings in a  
text format known as a *radiology report*. A key step in this  
diagnostic process is identifying the precise location and  
boundary of a *lesion*—an abnormal region with pathologi-  
cal changes [6]. This task is labor-intensive and demands  
substantial clinical expertise and analytical precision.

To alleviate physicians' workload in localizing pathologi-  
cal regions, there is a growing demand for automated lesion  
segmentation models in CXRs. Recently, vision-language  
models (VLMs) equipped with segmentation modules [22,  
23, 35] have emerged as a promising solution, as they can in-  
terpret diverse user-specific needs expressed through natural  
language instructions. However, despite the success of such  
VLMs in general-domain segmentation, their application to  
CXRs remains limited. Although prior studies [15, 25] have  
explored CXR lesion segmentation using text prompts, they

Table 1. Existing CXR datasets with spatial annotations for pathologic lesions.

Dataset	# Images	Spatial Annotation				Instruction-Answer Pair
		# Annotations	Type	Multi-Lesion	Method	
VinDr-CXR [33]	15K	9K	Bounding Box	✓	Manual	✗
Padchest-GR [8]	4.6K	7.7K	Bounding Box	✓	Manual	✗
MS-CXR [3]	1K	1.2K	Bounding Box	✓	Manual	✗
TBX-11K [28]	12K	1.2K	Bounding Box	✗	Manual	✗
SIIM-ACR [41]	13K	2.7K	Segmentation Mask	✗	Manual	✗
QaTa-COV19 [9]	121K	9.3K	Segmentation Mask	✗	Semi-Automated	✗
Danilov et al. [7]	1.4K	0.6K	Segmentation Mask	✓	Manual	✗
<b>MIMIC-ILS (Ours)</b>	192K	91K	Segmentation Mask	✓	Fully-Automated	✓

are limited to a single lesion type (e.g., COVID-19) and moreover require long, detailed expert-level medical descriptions based on tailored CXR review (e.g., “Bilateral pulmonary infection, two infected areas, upper right lung and upper left lung.”) as input. Such constraints make them impractical not only for physicians who aim to segment diverse lesion types across various anatomical subregions before closely reviewing the image themselves, but especially for non-experts who can hardly interpret CXR images at all.

To address these limitations, we propose a more user-friendly paradigm, namely *instruction-guided lesion segmentation* (ILS). In this paradigm, the model is required to process diverse user instructions, ranging from prompts that specify the lesion type and target location, to requests that look for abnormalities globally. If the requested lesion is not present, the model should reliably report its absence. Additionally, the model should be able to provide textual descriptions regarding a lesion’s location or type, even if not explicitly prompted by the user. However, a dataset to support such a versatile task has been unavailable, as constructing a suitable dataset for training and evaluation poses significant challenges—most notably the need for expert-curated mask annotations. Moreover, accurately pairing these masks with precise textual instructions in terms of anatomical locations and specific lesion types remains a highly complex task.

In this work, we introduce the first fully automated pipeline for constructing a large-scale ILS dataset for CXRs. The central challenge is: “How can we derive lesion masks and corresponding instruction-answer text pairs from raw images that contain no explicit annotations?” To address this, we leverage radiology reports as a key source of information for each image. Using paired image–report data, our two-stage pipeline integrates pre-trained vision models and large language models (LLMs) to extract high-confidence anomalous regions and structured textual information. By exploiting the consistency between these heterogeneous modalities, we generate high-quality lesion masks and diverse instruction–answer pairs. Applying our novel framework to MIMIC-CXR [19, 20]—a large, publicly available CXR–report dataset—we constructed **MIMIC-ILS**, a large-scale dataset consisting of 1.1M samples derived from 192K images and 91K lesion masks (Table 1).

Although several datasets [3, 7–9, 28, 33, 41] have tried to introduce spatial annotations in the CXR domain, they are unsuitable for direct use in our ILS paradigm (Table 1). Most provide only coarse bounding-box localization or single lesion type masks that are limited in scale due to reliance on expert annotations. Moreover, they also lack explicit links between mask annotations and textual instructions. MIMIC-ILS bridges these gaps by offering large-scale instruction–answer pairs, each paired with an auto-labeled segmentation mask and a detailed lesion profile. Despite being constructed entirely without human intervention, expert evaluations report a high acceptance rate of over 95% for this dataset.

Leveraging MIMIC-ILS, we train **ROSALIA** (Radiology Segmentation Assistant trained on a Lesion-grounded Instruction-Answer dataset), the first VLM designed for ILS in CXRs. Given user instructions, ROSALIA generates segmentation masks and textual descriptions (Fig. 1), handling a wide range of tasks, such as specific segmentation (e.g., “Segment the pneumonia in the right lung.”), generic segmentation (e.g., “Segment the opacity.”), and absence confirmation (e.g., “There is no atelectasis in the left lung base.”). This flexibility enables ROSALIA to effectively address diverse user needs, delivering tailored outputs for each request.

In summary, our contributions are threefold:

- We introduce a novel automated pipeline that generates lesion masks and corresponding instructions directly from CXRs without any human intervention. Using only image–report pairs, our method produces a large-scale dataset without requiring explicit manual processing.
- Applying our framework to MIMIC-CXR, we construct MIMIC-ILS, the first dataset for instruction-guided lesion segmentation (ILS) in CXRs. The resulting dataset is further validated by medical experts, confirming its high quality and the reliability of the construction process.
- To validate the utility of MIMIC-ILS, we introduce ROSALIA, the first VLM designed for ILS in CXRs. Trained on our million-scale dataset, ROSALIA interprets user instructions across diverse lesion types and locations, producing accurate lesion masks and descriptive outputs. As existing general and medical VLMs significantly struggle with this task, we will publicly release our dataset and model to support advances in fine-grained CXR lesion grounding.

130	<b>2. Related Work</b>		
131	<b>2.1. Lesion Segmentation and Datasets</b>		
132	Lesion segmentation aims to generate masks corresponding	instruction–answer pairs (Sec. 3.2).	181
133	to abnormal regions in medical images. Typically, models		
134	are trained on datasets where radiologists have directly annotated	<b>3.1. Grounded Lesion Mask Generation</b>	182
135	lesion masks. For CT and MRI, several studies [16, 17]	To construct our dataset, we use MIMIC-CXR [19, 20], a	183
136	have utilized public datasets that provide diverse tumor	large collection of CXR images paired with radiology reports.	184
137	masks [1, 2, 13]. In contrast, such pixel-level annotations	Each report is written by a radiologist and provides visual	185
138	are scarce in the CXR domain. While some datasets provide	descriptions of the corresponding CXR image. Based on this	186
139	only bounding boxes [8, 33], those that offer segmentation	dataset, we generate grounded lesion masks through four	187
140	masks usually focus on a single lesion type [9, 28, 41]. Conse-	sequential steps as illustrated in Fig. 2: (1) Report structuring	188
141	quently, existing models trained on these datasets are limited	and location mapping; (2) Spatial information extraction;	189
142	in their effective segmentation range for CXRs [42]. Our	(3) Lesion mask generation; (4) Location verification. The	190
143	work directly addresses this gap by constructing a compre-	details of each step are provided in Appendix A.	191
144	hensive, multi-type lesion segmentation dataset for CXRs.		
145	<b>2.2. Language-Guided Image Segmentation</b>	<b>Report Structuring and Location Mapping.</b> The first step	192
146	Language-guided image segmentation is the task of segment-	employs LLMs to convert radiology reports into a structured	193
147	ing a target specified by text. Early approaches to this task	form for later steps. Specifically, we instruct an LLM to	194
148	focused on aligning image features with text labels to gener-	transform each sentence describing an abnormal finding into	195
149	ate corresponding masks [24, 37, 40, 43]. More recently, ad-	a six-element tuple consisting of the following categories:	196
150	vancements in VLMs have enabled researchers to extend their	entity, sentence index, presence, certainty, location, and pre-	197
151	reasoning capabilities to segmentation [22, 23, 35]. These	dicted lesion type. The location element is then mapped to	198
152	models can generate an appropriate mask based on complex	one or more anatomical labels to ensure compatibility with	199
153	instructions that require real-world knowledge, such as “Seg-	the segmentation model used in subsequent processes. For	200
154	ment the object richest in vitamin C in this photo.”	example, if the second sentence in a given radiology report	201
155	Similar research has emerged in the medical domain, but	is “The lower lung opacity is pneumonia.”, its corresponding	202
156	current approaches remain limited. They usually rely on	output is (opacity, 2, positive, definitive, [right	203
157	simple prompts including class labels ( <i>e.g.</i> , “a computer-	lung base, left lung base], pneumonia). Here, the	204
158	ized tomography of a tumor”) [5, 27], which cannot handle	term “lower lung” in the original report is mapped to “right	205
159	sentence-level instructions. In the CXR domain specifically,	lung base” and “left lung base”.	206
160	recent VLMs have been trained using free-form text that de-	<b>Spatial Information Extraction.</b> The second step extracts	207
161	scribes the location and number of lesions [15, 25]. These	spatial information from CXRs using three distinct models:	208
162	approaches, however, expect users to have already reviewed	(1) RadEdit [34], a diffusion-based image editing model;	209
163	the CXR image, thus providing expert-level descriptions as	(2) CXAS [36], an anatomy segmentation model; and (3)	210
164	input. In contrast, our model allows users to obtain the lesion	a pretrained YOLO model for CXR lesion detection [32].	211
165	mask, its presence or absence, and type information even	These models are used respectively to generate an anomaly	212
166	without having to interpret the CXR images first.	map, anatomy masks, and lesion box masks, which serve	213
167	<b>3. Automatic Dataset Construction</b>	as visual cues for lesion mask generation in the subsequent	214
168	This section outlines our approach to automatically construct-	steps.	215
169	ing a large-scale dataset for training a model that generates	RadEdit takes an input image $x \in \mathbb{R}^{H \times W}$ containing	216
170	both lesion segmentation masks and corresponding textual	a lesion and the text prompt “No acute cardiopulmonary	217
171	descriptions in response to user instructions. The main chal-	process” and outputs an edited image $\hat{x}$ from which the lesion	218
172	lenges in this process are: (1) generating lesion masks directly	has been removed. We derive $x_{\text{ano}} \in [0, 1]^{H \times W}$ as:	219
173	from raw CXR images without explicit image annotations,		
174	(2) aligning appropriate instruction–answer texts with the	$x_{\text{ano}} = \frac{x - \hat{x}}{I_{\text{max}}},$	220
175	obtained masks, and (3) ensuring that the entire pipeline op-	where $I_{\text{max}}$ is the maximum possible pixel intensity ( <i>i.e.</i> , 255	221
176	erates in a fully automated, human-free manner. To address	for an 8-bit image). $x_{\text{ano}}$ provides morphological informa-	222
177	these challenges, our framework first extracts textual and	tion about hyperintense lesions, which are areas that appear	223
178	spatial information from image–report pairs and generates	brighter than the normal lung field. From $x_{\text{ano}}$ , we define	224
179	lesion masks followed by a verification process (Sec. 3.1).	anomaly map $\mathcal{A}$ as:	225
180	Using the verified lesion masks, we then construct diverse	$\mathcal{A} = \{(i, j) \mid (x_{\text{ano}})_{i,j} \geq \tau_{\text{ano}}\},$	226
		where $(i, j)$ represents a pixel coordinate and $\tau_{\text{ano}}$ is a thresh-	227
		old for anomaly pixels.	228

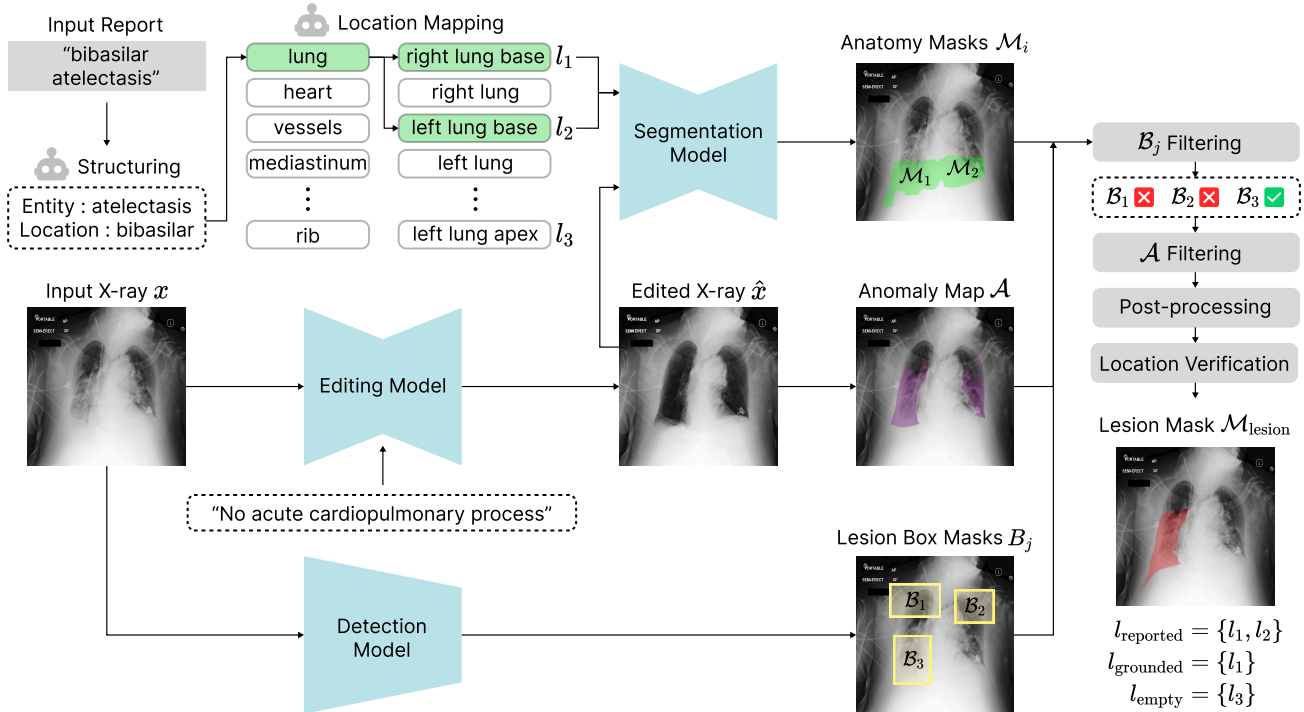


Figure 2. An overview of grounded lesion mask generation. **(Top-left)** Textual information is extracted from the radiology report during the report structuring and location mapping. **(Bottom-left and Center)** Pretrained vision models are also employed to produce spatial information. **(Right)** Finally, a lesion mask is generated by integrating this information. The verification step then confirms the grounded location ( $l_1$ ), identifies the empty location ( $l_3$ ) for negative sample generation, and discards the reported-but-ungrounded location ( $l_2$ ).

229 CXAS produces anatomy masks corresponding to each  
 230 location element in the previously derived structured  
 231 report tuples. We denote these masks as coordinate sets  $\{\mathcal{M}_i\}_{i=1}^n$ ,  
 232 where  $n$  is the number of anatomical labels mapped in the  
 233 previous step. Each  $\mathcal{M}_i$  contains the pixel coordinates for a  
 234 specific anatomy, serving as a spatial approximation of the  
 235 lesion location mentioned in the radiology report.

236 In parallel, the pretrained YOLO model is applied to  $x$  to  
 237 detect a diverse range of lesions. It outputs bounding boxes  
 238 that not only specify the locations of potential lesions but  
 239 also assign a confidence score to each detection. From these  
 240 results, we construct a set of lesion box masks,  $\{\mathcal{B}_j\}_{j=1}^m$ ,  
 241 where  $m$  denotes the number of detected boxes. Each  $\mathcal{B}_j$   
 242 represents the pixel coordinates enclosed by a bounding box,  
 243 accompanied by a confidence score  $conf_{\mathcal{B}_j} \in [0, 1]$ .

244 **Lesion Mask Generation.** With the three visual cues ex-  
 245 tracted from the previous step, the initial lesion masks can  
 246 be generated. Here, the anomaly map  $\mathcal{A}$  plays a central role,  
 247 representing a composite signal of all hyperintense lesions.  
 248 We decompose this signal into individual masks and align  
 249 them with the specific lesions described in the report. During  
 250 this process, the anatomy masks  $\{\mathcal{M}_i\}_{i=1}^n$ , lesion box masks  
 251  $\{\mathcal{B}_j\}_{j=1}^m$ , and the right and left lung masks ( $L_r$  and  $L_l$ )  
 252 are jointly used to select high-quality mask candidates.

253 The core of this filtering process, outlined in Algorithm 1,  
 254 selectively retains only appropriate candidates from the ini-  
 255 tially detected lesion box masks, based on four conditions:

( $c_1$ ) sufficient overlap with  $\{\mathcal{M}_i\}_{i=1}^n$ ; ( $c_2$ ) a high confidence  
 256 score; ( $c_3$ ) a high internal signal ratio from  $\mathcal{A}$  (*i.e.*, the ratio of  
 257 the intersection area between the box mask and  $\mathcal{A}$  to the area  
 258 of the box mask); and ( $c_4$ ) a sufficient size relative to either  
 259  $L_r$  or  $L_l$ . Conditions  $c_1$  and  $c_2$  ensure that the boxes align  
 260 with the reported locations and are likely to correspond to  
 261 true lesions. However, the  $\mathcal{A}$  can contain false negatives (*i.e.*,  
 262 coordinates that belong to actual lesion areas but are missing  
 263 from  $\mathcal{A}$ ), which may result in excessively small or even empty  
 264 masks. To mitigate this issue, conditions  $c_3$  and  $c_4$  are used  
 265 to retain only those boxes that contain strong lesion signals and  
 266 are large enough to allow meaningful segmentation. Once the  
 267 appropriate lesion box masks are selected, we extract from  
 268  $\mathcal{A}$  the connected components (*i.e.*, the individual, contig-  
 269 uous ‘islands’ in 2D space) that intersect with these selected  
 270 masks. This component then undergoes a post-processing  
 271 step involving small, noisy mask removal to produce the final,  
 272 refined lesion mask  $\mathcal{M}_{\text{lesion}}$  (see Appendix A.5 for further  
 273 details).  
 274

**Location Verification.** In the final step, we explicitly ver-  
 275 ify whether each lesion mask generated by Algorithm 1 has  
 276 been successfully grounded to the structured report. To assess  
 277 the grounding status, we define three types of locations: *re-*  
 278 *ported location*, *grounded location* and *empty location*. The  
 279 *reported location* is a set of anatomical labels extracted from  
 280 the previous location mapping with LLMs. Based on this set,  
 281 the *grounded location* is defined as a subset of the *reported*  
 282

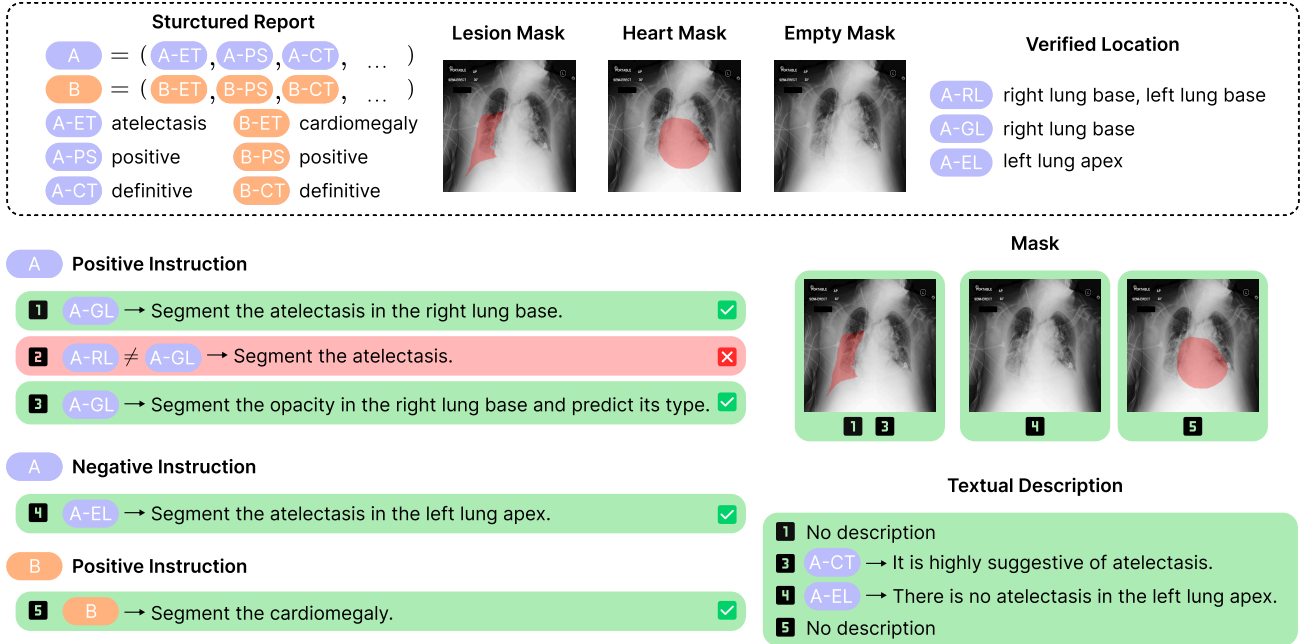


Figure 3. Instruction–answer pair generation process using the example report, “Bibasilar atelectasis. Cardiomegaly.” We utilize the elements extracted from the previous lesion mask generation process (see Fig. 2), indicated by the dashed box. Structured tuples (A&B in the top left) are converted to text instructions and mapped to their corresponding ground-truth masks and textual descriptions. Invalid instructions for lesions which lack a corresponding mask are excluded (colored as red), and only valid instructions are retained (colored as green). (ET: entity, PS: presence, CT: certainty, RL: reported location, GL: grounded location, EL: empty location)

### Algorithm 1: Lesion Mask Generation

**Input:** Anomaly map  $\mathcal{A}$ , anatomy masks  $\{\mathcal{M}_i\}_{i=1}^n$ , lesion box masks  $\{\mathcal{B}_j\}_{j=1}^m$  with confidences  $\{conf_{\mathcal{B}_j}\}_{j=1}^m$ , right lung mask  $L_r$ , left lung mask  $L_l$

**Output:** Final lesion mask  $\mathcal{M}_{\text{lesion}}$

```

1  $\mathcal{M}_{\text{lesion}} \leftarrow \emptyset$ ;
2  $\mathcal{M}_{\text{union}} \leftarrow \bigcup_{i=1}^n \mathcal{M}_i$ ;
3 foreach  $\mathcal{B}_j \in \{\mathcal{B}_j\}_{j=1}^m$  do
4    $c_1 \leftarrow \frac{|\mathcal{B}_j \cap \mathcal{M}_{\text{union}}|}{|\mathcal{B}_j \cup \mathcal{M}_{\text{union}}|} \geq \tau_{\text{anatomy}}$ ;
5    $c_2 \leftarrow conf_{\mathcal{B}_j} \geq \tau_{\text{conf}}$ ;
6    $c_3 \leftarrow \frac{|\mathcal{B}_j \cap \mathcal{A}|}{|\mathcal{B}_j|} \geq \tau_{\text{signal}}$ ;
7    $c_4 \leftarrow \left( \frac{|\mathcal{B}_j \cap L_r|}{|\mathcal{B}_j \cup L_r|} \geq \tau_{\text{size}} \right) \vee \left( \frac{|\mathcal{B}_j \cap L_l|}{|\mathcal{B}_j \cup L_l|} \geq \tau_{\text{size}} \right)$ ;
8   if  $c_1 \wedge c_2 \wedge c_3 \wedge c_4$  then
9      $C \leftarrow \text{FindIntersectingComponent}(\mathcal{B}_j, \mathcal{A})$ ;
10    if  $C$  is not empty then
11       $\mathcal{M}_{\text{new}} \leftarrow \text{Refine}(C)$ ;
12       $\mathcal{M}_{\text{lesion}} \leftarrow \mathcal{M}_{\text{lesion}} \cup \mathcal{M}_{\text{new}}$ ;
13 return  $\mathcal{M}_{\text{lesion}}$ 

```

283 location that spatially overlaps with a generated lesion mask,  
 284 confirming successful localization of the reported finding.  
 285 This location is derived from the anatomy masks  $\{\mathcal{M}_i\}_{i=1}^n$   
 286 that intersect with the selected lesion box masks during the  
 287 lesion mask generation. Finally, we introduce an *empty loca-*  
 288 *tion*, which refers to a lung region with no reported lesions  
 289 and is used to generate negative samples.

### 3.2. Instruction-Answer Pair Generation

290 With the information extracted from the previous process  
 291 (*i.e.*, grounded lesion mask generation), we build our dataset  
 292 for seven major lesion types found in CXRs: cardiomegaly,  
 293 pneumonia, atelectasis, opacity, consolidation, edema, and  
 294 effusion. These lesions are not only the most frequently  
 295 mentioned in radiology reports, but also clinically signif-  
 296 icant to be common annotation targets in other medical  
 297 datasets [3, 19, 20, 38]. For each lesion, we construct positive  
 298 instruction-answer pairs, which include a ground-truth lesion  
 299 mask. Negative pairs using an empty mask are also generated  
 300 to enable the model to confirm the absence of lesions. An  
 301 example of this pair generation process is shown in Fig. 3.  
 302 Please refer to Appendix B and C for the lesion descriptions  
 303 and specific dataset generation process.  
 304

**Instruction Types and Limitations.** We consider three types  
 305 of segmentation instructions (Table 2). A *basic* instruction  
 306 specifies both the segmentation target and its location. The  
 307 location can be a broad region (such as left lung or right  
 308 lung), one of eight more specific zones (apical, upper, mid,  
 309 and lower zones for each lung), or a combination of these  
 310 regions. In contrast, a *global* instruction specifies only the  
 311 segmentation target. A *lesion inference* instruction asks the  
 312 model to predict the type of lesion represented by an opacity  
 313 within a given location. The generation of these instructions  
 314 is inherently constrained by the grounded lesion mask gener-  
 315 ation. For example, a global instruction becomes invalid if  
 316

317 the generated mask captures only part of the lesion. To ad-  
 318 dress this, our framework dynamically produces only those  
 319 instruction–answer pairs that are valid given the grounding  
 320 information available for each image.

Table 2. Templates for each question type. Each type includes answer templates for both positive and negative cases, with the negative answers positioned in the last row of each cell.

Type	Role	Template
Basic	Instruction	Segment the [Target] in the [Location].
	Answer	[SEG] It is located in the [Location]. [SEG] There is no [Target] in the [Location].
Global	Instruction	Segment the [Target].
	Answer	[SEG] It is located in the [Location]. [SEG] There is no [Target].
Lesion Inference	Instruction	Segment the opacity in the [Location] and predict its type.
	Answer	[SEG] It is highly suggestive of [Lesion]. [SEG] It possibly reflects [Lesion]. [SEG] There is no opacity in the [Location].

321 **Instruction Generation.** The instruction generation process  
 322 begins by creating a basic instruction for each grounded lesion.  
 323 Next, we determine whether a global instruction can  
 324 be generated. The global instruction is created only when  
 325 the *grounded location* and the *reported location* are identical.  
 326 Separately, we generate lesion inference instructions  
 327 by transforming the basic instructions for *pneumonia*, *at-*  
 328 *electasis*, and *edema*, replacing these specific lesion types  
 329 with *opacity*. This transformation is motivated by the fact  
 330 that these findings are all specific types of “opacity,” a more  
 331 fundamental visual concept in medical imaging. Negative  
 332 samples are generated by (1) selecting lesion types that are  
 333 not mentioned or explicitly negated in the radiology report; or  
 334 (2) utilizing *empty locations* to substitute the original location  
 335 in the basic instruction of a positive sample.

336 **Answer Generation.** Each answer consists of a lesion mask  
 337 and a textual description. The answer lesion masks for posi-  
 338 tive pairs are determined differently depending on whether  
 339 they are organ-level or localized abnormalities. For cardiomegaly,  
 340 we utilize a heart mask as its corresponding lesion mask since  
 341 this condition is defined by the state of a specific organ [11].  
 342 In contrast, localized abnormalities (*e.g.*, pneumonia or effusion)  
 343 can appear in variable locations, so for these findings, we use  
 344 the lesion masks generated in Sec. 3.1. For negative pairs, an  
 345 empty mask is used. As for the textual description, it is also  
 346 provided for both positive and negative samples. Specifically,  
 347 the answer template for lesion inference incorporates a certainty  
 348 level.

## 349 4. MIMIC-ILS Dataset

350 Our final dataset, MIMIC-ILS, consists of 1.1M instruction-  
 351 answer pairs (135K positive and 930K negative samples)  
 352 derived from 192K MIMIC-CXR images. This final image

set is obtained by first filtering out low-quality images (*e.g.*,  
 images with extreme contrast issues), and then excluding any  
 images for which no instruction–answer pairs are generated  
 through our pipeline in Sec. 3. The positive samples are gen-  
 erated from 91K unique lesion masks, where each mask can  
 be associated with multiple instruction–answer samples. The  
 resulting dataset covers seven distinct lesion types, and the  
 overall statistics are presented in Fig. 4. Following the official  
 MIMIC-CXR split, the dataset is divided into 1M training  
 samples, 8.2K validation samples, and 12K test samples. De-  
 tails on quality control and a distribution of MIMIC-ILS are  
 presented in Appendix D and E.

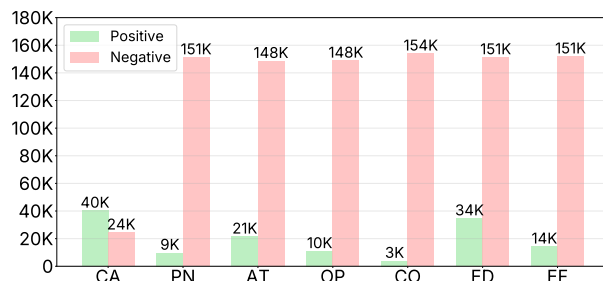


Figure 4. Distribution of MIMIC-ILS dataset. The y-axis indicates the number of samples, and the x-axis represents the lesion type. (CA: cardiomegaly, PN: pneumonia, AT: atelectasis, OP: opacity, CO: consolidation, ED: edema, EF: effusion)

**Human Evaluation.** To assess the quality of MIMIC-ILS,  
 an expert review was conducted by four radiation oncologists  
 specializing in lesion contouring on medical images. For  
 the test set samples, clinicians classified each case as either  
 acceptable or unacceptable based on mask quality. Positive  
 cases were reviewed by all experts, while negatives were split  
 among them. Any sample judged unacceptable by at least  
 one expert was excluded from the final test set, and the results  
 are summarized in Table 3. Among the 10.5K mask samples  
 initially reviewed, 96.3% were rated as acceptable and finally  
 included in the test set. More details on the expert profiles  
 and quality assessment are provided in Appendix E.

Table 3. Acceptance rate and number of evaluated samples for the human evaluation. Each sample corresponds to a unique combination of lesion mask, target, and location.

Expert	Total		Positive		Negative	
	Rate (%)	# Samples	Rate (%)	# Samples	Rate (%)	# Samples
Expert A	96.1	4,055	95.6	1,841	96.5	2,214
Expert B	97.2	3,968	96.0	1,841	98.3	2,217
Expert C	98.7	4,002	99.8	1,841	97.8	2,161
Expert D	97.6	4,039	96.9	1,841	98.2	2,198
Overall	96.3	10,541	90.1	1,841	97.7	8,700

## 377 5. Model Training

Using MIMIC-ILS, we train our ILS model, ROSALIA. The  
 model adopts the architecture of LISA [22], which demon-  
 strated strong zero-shot language-guided segmentation per-  
 formance in the general domain. As illustrated in Fig. 5, the

382 architecture integrates a VLM backbone with the Segment  
 383 Anything Model (SAM) [21]. The VLM processes both the  
 384 image and the input instruction to produce a special token,  
 385 [SEG], along with its textual description. This [SEG]  
 386 embedding is then passed to SAM together with the input  
 387 image for mask prediction. Within SAM, the frozen image  
 388 encoder extracts embeddings from the image, and the mask  
 389 decoder integrates these embeddings with the hidden embed-  
 390 ding of [SEG] token to generate the final mask.

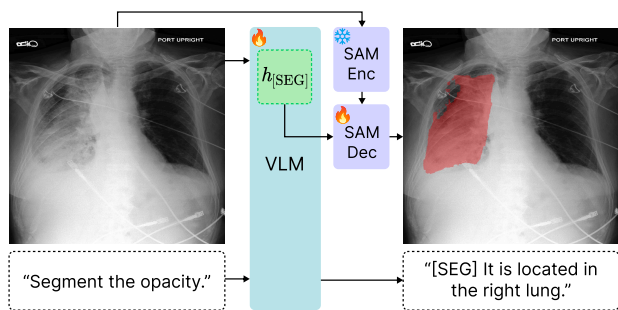


Figure 5. Overview of ROSALIA. The architecture integrates a VLM with the SAM. The VLM takes a CXR image and a segmentation instruction as input, generating both a textual description and a special [SEG] token. The hidden embedding of this [SEG] token is then passed to SAM’s decoder to produce the final mask.

391 The overall loss function  $\mathcal{L}$  consists of two components:  
 392 (1) a language loss and (2) a mask loss. It is formulated as:

$$\mathcal{L} = \lambda_{\text{txt}}\mathcal{L}_{\text{txt}} + \mathcal{L}_{\text{mask}},$$

$$\mathcal{L}_{\text{mask}} = \lambda_{\text{bce}}\mathcal{L}_{\text{bce}} + \lambda_{\text{dice}}\mathcal{L}_{\text{dice}}.$$

395  $\mathcal{L}_{\text{txt}}$  denotes the autoregressive cross-entropy loss for the  
 396 answer text, and  $\mathcal{L}_{\text{mask}}$  represents the segmentation loss com-  
 397 puted between the ground-truth mask and the predicted fore-  
 398 ground probability map, which combines the binary cross-  
 399 entropy loss  $\mathcal{L}_{\text{bce}}$  and the DICE [30] loss  $\mathcal{L}_{\text{dice}}$ . The  $\lambda_{\text{txt}}$ ,  $\lambda_{\text{bce}}$ ,  
 400 and  $\lambda_{\text{dice}}$  are coefficients for each loss term.

## 401 6. Experiments

### 402 6.1. Implementation Details

403 **Training Details.** ROSALIA is built on the LISA-7B archi-  
 404 tecture and is fine-tuned from its original checkpoint [22].  
 405 Following LISA, we adopted LLaVA [26] as the VLM back-  
 406 bone and employed the largest version of SAM (SAM-H).  
 407 LoRA [14] fine-tuning was applied to the VLM with a rank  
 408 of 128 and an alpha of 256, while the mask decoder was fully  
 409 fine-tuned. The epochs and the initial learning rate were set to  
 410 15 and 0.0003, respectively, using the AdamW optimizer [29].  
 411 The total batch size was 256, and the ratio of positive to neg-  
 412 ative samples was maintained at 1:1 in each mini-batch. The  
 413 loss coefficients  $\lambda_{\text{txt}}$ ,  $\lambda_{\text{bce}}$ , and  $\lambda_{\text{dice}}$  were set to 0.5, 5, and  
 414 1, respectively, and the DICE loss was computed only for  
 415 positive samples. Further model training details are described  
 416 in the Appendix F.

**Baseline Models.** Since we present the first dataset for ILS  
 in CXRs, no existing models have been directly trained on  
 our proposed task. Nonetheless, we evaluated several models  
 from both the general domain (LISA [22], Text4Seg [23],  
 PixelLM [35]) and the medical domain (BiomedParse [42],  
 RecLMIS [15], IMIS-Net [5]), which can take an image and  
 text as input to produce a segmentation output.

**Evaluation Metrics.** We used three metrics to evalu-  
 ate model performance. For positive samples, we used  
 Intersection-over-Union (IoU)-based measures: gIoU and  
 cIoU [22]. gIoU is the average IoU across samples, while  
 cIoU is the ratio of total intersection to total union across the  
 dataset. For negative cases, we used empty-target accuracy  
 (N-Acc.), the proportion of samples correctly predicted to  
 have no masks [39].

### 432 6.2. Main Results

433 Table 4 presents the results of the baselines and our proposed  
 434 model on the MIMIC-ILS test set. While existing VLM-based  
 435 segmentation models from both the general and medical do-  
 436 mains struggle with the ILS task, ROSALIA achieves notably  
 437 high performance. In particular, not only do these baselines  
 438 yield low IoU scores on positive cases, but they also fre-  
 439 quently fail on empty-target cases, where the N-Acc. rate  
 440 is nearly zero in most instances. These results highlight the  
 441 need for a dedicated dataset to effectively address the ILS  
 442 task in CXRs. Furthermore, the strong results of ROSALIA  
 443 on the physician-verified test set demonstrate that the training  
 444 set of MIMIC-ILS serves as a high-quality resource—even  
 445 without manual expert filtering.

Table 4. Segmentation results (%) on the MIMIC-ILS test set. “N-Acc.” denotes the accuracy of correctly predicting empty targets. ¶ indicates medical domain baselines. The best and second-best results are marked in **bold** and underline, respectively.

Model	gIoU	cIoU	N-Acc.
LISA-7B [22]	8.3	12.8	0.7
LISA-13B [22]	8.9	12.2	0.0
Text4Seg [23]	6.1	10.3	20.6
PixelLM-7B [35]	9.2	11.8	0.0
PixelLM-13B [35]	12.8	15.4	0.0
BiomedParse ¶ [42]	<b>23.8</b>	18.5	0.6
RecLMIS ¶ [15]	21.9	<u>18.6</u>	0.0
IMIS-Net ¶ [5]	9.8	11.8	<u>21.6</u>
<b>ROSALIA (Ours)</b>	<b>71.2</b>	<b>75.6</b>	<b>91.8</b>

446 Table 5 presents the performance of ROSALIA across  
 447 different lesion types. The overall gIoU exceeds 0.7, indi-  
 448 cating that more than 80% of the regions overlap between  
 449 the predicted and ground-truth masks when the two are of  
 450 similar size. Even for the lesion type with the lowest gIoU,  
 451 the score remains above 0.55, suggesting over 70% regional  
 452 overlap under similar mask sizes between the ground truth  
 453 and predictions.

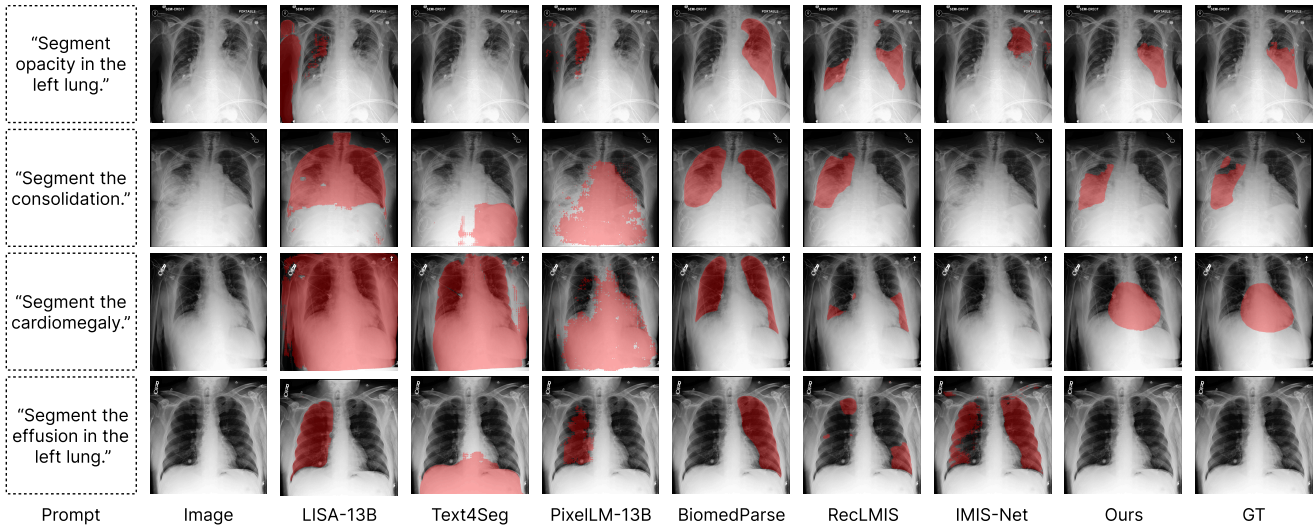


Figure 6. Visualized inference results of ROSALIA and baseline models. The first three rows show results for positive cases, while the last row presents results for negative cases with an empty target mask. Additional examples are demonstrated in Appendix G.

Table 5. Segmentation performance (%) of ROSALIA for each lesion type.

Lesion	gIoU	cIoU	N-Acc.
Cardiomegaly	89.0	89.0	85.8
Pneumonia	57.2	60.4	97.1
Atelectasis	60.2	58.7	91.7
Opacity	60.5	64.2	85.0
Consolidation	61.9	65.6	91.2
Edema	64.8	66.6	92.2
Effusion	60.3	59.6	90.4
<b>Total</b>	<b>71.2</b>	<b>75.6</b>	<b>91.8</b>

We also evaluate the accuracy of text responses across different question types, as shown in Table 6. A response is considered correct only when both the template and all variables for each question type (denoted by square brackets in Table 2) exactly match the structured ground-truth information. Despite this strict criterion, ROSALIA achieves high accuracy across most question types (see Appendix G for text accuracy of each lesion type).

Table 6. Text response accuracy (%) of ROSALIA.

Type	Overall	Basic	Global	Lesion Inf.
Positive	90.7	95.4	93.7	75.1
Negative	95.3	96.9	82.3	90.6
<b>Total</b>	<b>94.4</b>	<b>96.8</b>	<b>88.8</b>	<b>84.8</b>

### 6.3. Qualitative Results

Fig. 6 presents qualitative examples from each model for the ILS task. The baseline models largely fail, either producing entirely incorrect masks or segmenting the whole anatomical regions (e.g., the left or right lung). In contrast, ROSALIA accurately segments only the lesion specified in the instruction within the designated region and correctly identifies

empty-target cases. Additionally, Fig. 7 demonstrates the outputs produced from diverse instructions applied to the same input image. Although multiple lesions coexist in the image, ROSALIA accurately interprets each instruction and generates results tailored to the user’s specific request. This highlights the model’s ability to handle diverse lesion types and locations of interest.

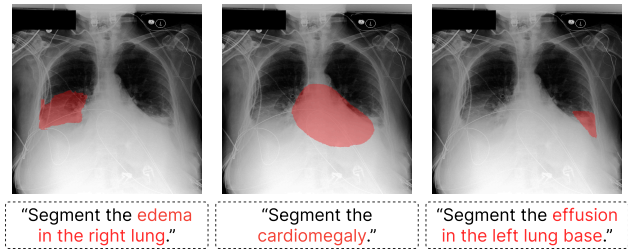


Figure 7. Examples of outputs from different instructions applied to the same image. Among the multiple lesions present, ROSALIA can selectively segment only the lesion and location of interest.

### 7. Conclusion

In this study, we introduce MIMIC-ILS, the first dataset for instruction-guided lesion segmentation in CXRs, along with ROSALIA, a VLM developed for this new paradigm. Our automated pipeline enables the construction of this million-scale dataset, and expert evaluations show a remarkably high acceptance rate, confirming the quality and reliability of our fully human-free data generation process. Trained on MIMIC-ILS, ROSALIA demonstrates a comprehensive ability to generate accurate lesion segmentations and textual responses across diverse user instructions. These findings indicate that MIMIC-ILS and ROSALIA offer a strong foundation for advancing research on fine-grained lesion grounding in the CXR domain.

490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545

## References

- [1] Michela Antonelli, Annika Reinke, Spyridon Bakas, Keyvan Farahani, Annette Kopp-Schneider, Bennett A Landman, Geert Litjens, Bjoern Menze, Olaf Ronneberger, Ronald M Summers, et al. The medical segmentation decathlon. *Nature communications*, 13(1):4128, 2022. 3
- [2] Patrick Bilic, Patrick Christ, Hongwei Bran Li, Eugene Vorontsov, Avi Ben-Cohen, Georgios Kaissis, Adi Szeskin, Colin Jacobs, Gabriel Efrain Humpire Mamani, Gabriel Chartrand, et al. The liver tumor segmentation benchmark (lits). *Medical image analysis*, 84:102680, 2023. 3
- [3] Benedikt Boecking, Naoto Usuyama, Shruthi Bannur, Daniel C Castro, Anton Schwaighofer, Stephanie Hyland, Maria Wetscherek, Tristan Naumann, Aditya Nori, Javier Alvarez-Valle, et al. Making the most of text semantics to improve biomedical vision–language processing. In *European conference on computer vision*, pages 1–21. Springer, 2022. 2, 5
- [4] Joshua Broder. Imaging the chest: the chest radiograph. *Diagnostic imaging for the emergency physician*, page 185, 2011. 1
- [5] Junlong Cheng, Bin Fu, Jin Ye, Guoan Wang, Tianbin Li, Haoyu Wang, Ruoyu Li, He Yao, Junren Cheng, JingWen Li, et al. Interactive medical image segmentation: A benchmark dataset and baseline. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 20841–20851, 2025. 3, 7
- [6] Sherri de Coronado, Margaret W Haber, Nicholas Sioutos, Mark S Tuttle, and Lawrence W Wright. Nci thesaurus: using science-based terminology to integrate cancer research results. In *MEDINFO 2004*, pages 33–37. IOS Press, 2004. 1
- [7] Viacheslav V Danilov, Alex Proutski, Alex Karpovsky, Alexander Kirpich, Diana Litmanovich, Dato Nefaridze, Oleg Talalov, Semyon Semyonov, Vladimir Koniukhovskii, Vladimir Shvartc, et al. Indirect supervision applied to covid-19 and pneumonia classification. *Informatics in Medicine Unlocked*, 28:100835, 2022. 2
- [8] Daniel Coelho de Castro, Aurelia Bustos, Shruthi Bannur, Stephanie L Hyland, Kenza Bouzid, Maria Teodora Wetscherek, Maria Dolores Sánchez-Valverde, Lara Jaques-Pérez, Lourdes Pérez-Rodríguez, Kenji Takeda, et al. Padchest-gr: A bilingual chest x-ray dataset for grounded radiology report generation. *NEJM AI*, 2(7):AIdbp2401120, 2025. 2, 3
- [9] Aysen Degerli, Serkan Kiranyaz, Muhammad EH Chowdhury, and Moncef Gabbouj. Osegnet: Operational segmentation network for covid-19 detection using chest x-ray images. In *2022 IEEE International Conference on Image Processing (ICIP)*, pages 2306–2310. IEEE, 2022. 2, 3
- [10] Nicolas Gaggion, Candelaria Mosquera, Martina Aineseder, Lucas Mansilla, Diego Milone, and Enzo Ferrante. Chexmask database: a large-scale dataset of anatomical segmentation masks for chest x-ray images. 1
- [11] Nicolas Gaggion, Lucas Mansilla, Candelaria Mosquera, Diego H. Milone, and Enzo Ferrante. Improving anatomical plausibility in medical image segmentation via hybrid graph neural networks: applications to chest x-ray analysis. *IEEE Transactions on Medical Imaging*, 2022. 6
- [12] Nicolás Gaggion, Candelaria Mosquera, Lucas Mansilla, Julia Mariel Saidman, Martina Aineseder, Diego H Milone, and Enzo Ferrante. Chexmask: a large-scale dataset of anatomical segmentation masks for multi-center chest x-ray images. *Scientific Data*, 11(1):511, 2024. 1
- [13] Nicholas Heller, Sean McSweeney, Matthew Thomas Peterson, Sarah Peterson, Jack Rickman, Bethany Stai, Resha Tejpaul, Makinna Oestreich, Paul Blake, Joel Rosenberg, et al. An international challenge to use artificial intelligence to define the state-of-the-art in kidney and kidney tumor segmentation in ct imaging., 2020. 3
- [14] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. 7
- [15] Xiaoshuang Huang, Hongxiang Li, Meng Cao, Long Chen, Chenyu You, and Dong An. Cross-modal conditioned reconstruction for language-guided medical image segmentation. *IEEE Transactions on Medical Imaging*, 2024. 1, 3, 7
- [16] Yankai Jiang, Zhongzhen Huang, Rongzhao Zhang, Xiaofan Zhang, and Shaoting Zhang. Zept: Zero-shot pan-tumor segmentation via query-disentangling and self-prompting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11386–11397, 2024. 3
- [17] Yankai Jiang, Wenhui Lei, Xiaofan Zhang, and Shaoting Zhang. Unleashing the potential of vision-language pre-training for 3d zero-shot lesion segmentation via mask-attribute alignment. *arXiv preprint arXiv:2410.15744*, 2024. 3
- [18] Alistair Johnson, Matt Lungren, Yifan Peng, Zhiyong Lu, Roger Mark, Seth Berkowitz, and Steven Horng. Mimic-cxr-jpg-chest radiographs with structured labels. *PhysioNet*, 101: 215–220, 2019. 1
- [19] A Johnson, T Pollard, R Mark, S Berkowitz, and S Horng. Mimic-cxr database (version 2.1. 0). physionet. rrid: Scr\_007345, 2024. 2, 3, 5, 1
- [20] Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Roger G Mark, and Steven Horng. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data*, 6(1):317, 2019. 2, 3, 5, 1
- [21] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. 7
- [22] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9579–9589, 2024. 1, 3, 6, 7
- [23] Mengcheng Lan, Chaofeng Chen, Yue Zhou, Jiaying Xu, Yiping Ke, Xinjiang Wang, Litong Feng, and Wayne Zhang. Text4seg: Reimagining image segmentation as text generation. *arXiv preprint arXiv:2410.09855*, 2024. 1, 3, 7

- 604 [24] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen  
605 Koltun, and René Ranftl. Language-driven semantic seg-  
606 mentation. *arXiv preprint arXiv:2201.03546*, 2022. 3
- 607 [25] Zihan Li, Yunxiang Li, Qingde Li, Puyang Wang, Dazhou  
608 Guo, Le Lu, Dakai Jin, You Zhang, and Qingqi Hong. Lvit:  
609 language meets vision transformer in medical image segmen-  
610 tation. *IEEE transactions on medical imaging*, 43(1):96–107,  
611 2023. 1, 3
- 612 [26] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee.  
613 Visual instruction tuning. *Advances in neural information  
614 processing systems*, 36:34892–34916, 2023. 7
- 615 [27] Jie Liu, Yixiao Zhang, Jie-Neng Chen, Junfei Xiao, Yongyi  
616 Lu, Bennett A Landman, Yixuan Yuan, Alan Yuille, Yucheng  
617 Tang, and Zongwei Zhou. Clip-driven universal model for  
618 organ segmentation and tumor detection. In *Proceedings of  
619 the IEEE/CVF international conference on computer vision*,  
620 pages 21152–21164, 2023. 3
- 621 [28] Yun Liu, Yu-Huan Wu, Yunfeng Ban, Huifang Wang, and  
622 Ming-Ming Cheng. Rethinking computer-aided tuberculosis  
623 diagnosis. In *Proceedings of the IEEE/CVF conference on  
624 computer vision and pattern recognition*, pages 2646–2655,  
625 2020. 2, 3
- 626 [29] Ilya Loshchilov and Frank Hutter. Decoupled weight decay  
627 regularization. *arXiv preprint arXiv:1711.05101*, 2017. 7
- 628 [30] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi.  
629 V-net: Fully convolutional neural networks for volumetric  
630 medical image segmentation. In *2016 fourth international  
631 conference on 3D vision (3DV)*, pages 565–571. Ieee, 2016. 7
- 632 [31] Duc Nguyen, DungNB, Ha Q. Nguyen, Julia Elliott, Nguyen-  
633 ThanhNhan, and Phil Culliton. Vinbigdata chest x-ray abnor-  
634 malities detection. [https://kaggle.com/competitions/  
635 vinbigdata-chest-xray-abnormalities-detection](https://kaggle.com/competitions/vinbigdata-chest-xray-abnormalities-detection),  
636 2020. Kaggle. 1
- 637 [32] Dung Nguyen, Minh Khoi Ho, Huy Ta, Thanh Tam Nguyen,  
638 Qi Chen, Kumar Rav, Quy Duong Dang, Satwik Ramchandre,  
639 Son Lam Phung, Zhibin Liao, Minh-Son To, Johan Verjans,  
640 Phi Le Nguyen, and Vu Minh Hieu Phan. Localizing be-  
641 fore answering: A benchmark for grounded medical visual  
642 question answering. In *Proceedings of the Thirty-Fourth  
643 International Joint Conference on Artificial Intelligence (IJ-  
644 CAI 2025)*, pages 7670–7676, 2025. 3
- 645 [33] Ha Q Nguyen, Khanh Lam, Linh T Le, Hieu H Pham, Dat Q  
646 Tran, Dung B Nguyen, Dung D Le, Chi M Pham, Hang TT  
647 Tong, Diep H Dinh, et al. Vindr-cxr: An open dataset of chest  
648 x-rays with radiologist’s annotations. *Scientific Data*, 9(1):  
649 429, 2022. 2, 3, 1
- 650 [34] Fernando Pérez-García, Sam Bond-Taylor, Pedro P Sanchez,  
651 Boris van Breugel, Daniel C Castro, Harshita Sharma,  
652 Valentina Salvatelli, Maria TA Wetscherek, Hannah Richard-  
653 son, Matthew P Lungren, et al. Radedit: stress-testing biomed-  
654 ical vision models via diffusion image editing. In *European  
655 Conference on Computer Vision*, pages 358–376. Springer,  
656 2024. 3, 1
- 657 [35] Zhongwei Ren, Zhicheng Huang, Yunchao Wei, Yao Zhao,  
658 Dongmei Fu, Jiashi Feng, and Xiaojie Jin. Pixellm: Pixel  
659 reasoning with large multimodal model. In *Proceedings of  
660 the IEEE/CVF Conference on Computer Vision and Pattern  
661 Recognition*, pages 26374–26383, 2024. 1, 3, 7
- [36] Constantin Seibold, Alexander Jaus, Matthias A Fink, Moon  
662 Kim, Simon Reiß, Ken Herrmann, Jens Kleesiek, and Rainer  
663 Stiefelhagen. Accurate fine-grained segmentation of human  
664 anatomy in radiographs via volumetric pseudo-labeling. *arXiv  
665 preprint arXiv:2306.03934*, 2023. 3, 1
- [37] Gyungin Shin, Weidi Xie, and Samuel Albanie. Reco: Re-  
666 trieve and co-segment for zero-shot transfer. *Advances in  
667 Neural Information Processing Systems*, 35:33754–33767,  
668 2022. 3
- [38] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Moham-  
669 madhadi Bagheri, and Ronald M Summers. Chestx-ray8:  
670 Hospital-scale chest x-ray database and benchmarks on  
671 weakly-supervised classification and localization of common  
672 thorax diseases. In *Proceedings of the IEEE conference on  
673 computer vision and pattern recognition*, pages 2097–2106,  
674 2017. 5
- [39] Zhuofan Xia, Dongchen Han, Yizeng Han, Xuran Pan, Shiji  
675 Song, and Gao Huang. Gsva: Generalized segmentation  
676 via multimodal large language models. In *Proceedings of  
677 the IEEE/CVF Conference on Computer Vision and Pattern  
678 Recognition*, pages 3858–3869, 2024. 7
- [40] Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon,  
679 Thomas Breuel, Jan Kautz, and Xiaolong Wang. Groupvit:  
680 Semantic segmentation emerges from text supervision. In  
681 *Proceedings of the IEEE/CVF conference on computer vision  
682 and pattern recognition*, pages 18134–18144, 2022. 3
- [41] Anna Zawacki, Carol Wu, George Shih, Julia Elliott, Mikhail  
683 Fomitchev, Mohannad Hussain, ParasLakhani, Phil Cullit-  
684 on, and Shunxing Bao. Siim-acr pneumothorax segmen-  
685 tation. [https://kaggle.com/competitions/siim-acr-  
686 pneumothorax-segmentation](https://kaggle.com/competitions/siim-acr-pneumothorax-segmentation), 2019. Kaggle. 2, 3
- [42] Theodore Zhao, Yu Gu, Jianwei Yang, Naoto Usuyama,  
687 Ho Hin Lee, Tristan Naumann, Jianfeng Gao, Angela Crab-  
688 tree, Jacob Abel, Christine Mounq-Wen, et al. Biomedparse: a  
689 biomedical foundation model for image parsing of everything  
690 everywhere all at once. *arXiv preprint arXiv:2405.12971*,  
691 2024. 3, 7
- [43] Chong Zhou, Chen Change Loy, and Bo Dai. Extract free  
692 dense labels from clip. In *European conference on computer  
693 vision*, pages 696–712. Springer, 2022. 3