

# MeDSLIP: Medical Dual-Stream Language-Image Pre-training with Pathology-Anatomy Semantic Alignment

Wenrui Fan, Mohammad N.I. Suvon, Shuo Zhou, Xianyuan Liu, Samer Alabed, Venet Osmani, Andrew J. Swift, Chen Chen, and Haiping Lu, *Senior Member, IEEE*

**Abstract**—Pathology and anatomy are two essential groups of semantics in medical data. Pathology describes what the diseases are, while anatomy explains where the diseases occur. They describe diseases from different perspectives, providing complementary insights into diseases. Thus, properly understanding these semantics and their relationships can enhance medical vision-language models (VLMs). However, pathology and anatomy semantics are usually entangled in medical data, hindering VLMs from explicitly modeling these semantics and their relationships. To address this challenge, we propose MeDSLIP, a novel Medical Dual-Stream Language-Image Pre-training pipeline, to disentangle pathology and anatomy semantics and model the relationships between them. We introduce a dual-stream mechanism in MeDSLIP to explicitly disentangle medical semantics into *pathology-relevant* and *anatomy-relevant* streams and align visual and textual information within each stream. Furthermore, we propose an interaction modeling module with *prototypical contrastive learning loss* and *intra-image contrastive learning loss* to regularize the relationships between pathology and anatomy semantics. We apply MeDSLIP to chest X-ray analysis and conduct comprehensive evaluations with four benchmark datasets: NIH CXR14, RSNA Pneumonia, SIIM-ACR Pneumothorax, and COVIDx CXR-4. The results demonstrate MeDSLIP's superior generalizability and transferability across different scenarios. The code is available at <https://github.com/Shef-AIRE/MeDSLIP>, and the pre-trained model is released at <https://huggingface.co/pykale/MeDSLIP>.

**Index Terms**—Chest X-ray, Medical Vision-Language Model, Semantic Alignment.

## I. INTRODUCTION

**P**ATHOLOGY and anatomy are two essential groups of semantics in medical images and associated reports. Pathology focuses on the nature and characteristics of diseases, explaining what the abnormalities are. Anatomy, on the other hand, provides the structural and locational context, describing where these abnormalities occur [1]. For instance, in the sentence, “Opacity is observed on the bilateral lungs, and deformity of posterior ribs is noted,” “opacity” and “deformity” are pathology semantics, while “lungs” and “ribs” are anatomy semantics [2].

Pathology and anatomy semantics describe diseases from different perspectives, offering complementary insights into understanding diseases [3], [4]. Therefore, clearly modeling these semantics and their relationships can significantly enhance disease understanding via medical vision-language models (VLMs) [5], [6], thereby improving performance on key medical tasks such as medical image analysis [7], [8].

However, pathology and anatomy semantics are deeply entangled in medical contexts: Pathology semantics are often contextualized within specific anatomical regions. This entanglement hinders medical VLMs from explicitly understanding pathology and anatomy semantics and their intrinsic relationships, which further leads to the underutilization of the information in data. Hence, it will be beneficial to disentangle pathology and anatomy semantics from data and model their relationship with explicit guidance.

Most existing medical VLMs for extracting semantics from medical data are text-centric. They often prioritize textual hierarchy [9]–[19] and semantics [20]–[23] in medical reports with limited exploitation of the intricate visual semantics presented in medical images. This imbalance leads to underutilization of pathology and anatomy semantics in medical images, missing the opportunity to use their complementary insights for more effective image analysis.

Two VLM pre-training methods to align visual and textual information are *hierarchical alignment* and *semantic alignment*. Both are text-centric. (1) Hierarchical alignment stratifies textual information into different levels and aligns

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

Wenrui Fan, Mohammad N.I. Suvon, Shuo Zhou, Xianyuan Liu, and Haiping Lu are with Centre for Machine Intelligence and School of Computer Science, University of Sheffield, S1 4DP Sheffield, U.K. (e-mail: wenrui.fan@sheffield.ac.uk; m.suvon@sheffield.ac.uk; shuo.zhou@sheffield.ac.uk; xianyuan.liu@sheffield.ac.uk; h.lu@sheffield.ac.uk) (Corresponding author: Haiping Lu).

Samer Alabed and Andrew J. Swift are with School of Medicine and Population Health, and INSIGNEO, Institute for in Silico Medicine, University of Sheffield, S10 2TN Sheffield, and Department of Clinical Radiology, Sheffield Teaching Hospitals, S10 2JF Sheffield, U.K. (e-mail: s.alabed@sheffield.ac.uk; a.j.swift@sheffield.ac.uk).

Venet Osmani is with Digital Environment Research Institute, Queen Mary University of London, E1 1HH London, U.K. and School of Computer Science, University of Sheffield, S1 4DP Sheffield, U.K. (e-mail: v.osmani@qmul.ac.uk).

Chen Chen is with School of Computer Science, University of Sheffield, S1 4DP Sheffield, U.K. and Department of Computing, Imperial College London, SW7 2AZ London, U.K. (e-mail: chen.chen2@sheffield.ac.uk).

visual features accordingly [9]–[19]. This text stratification is usually based on textual hierarchies of the medical reports, such as syntax [17] or discourse [18] hierarchies. (2) Semantic alignment focuses on semantic concepts in the data but primarily on textual semantics [20]–[24]. It usually extracts key semantics from raw medical reports and enhances those textual semantics with prior knowledge from humans (e.g., domain knowledge [20], [23], [24], knowledge graph [21], [22], etc).

Emphasizing texts, most current medical VLMs learn visual semantics automatically in pre-training without explicit guidance, leaving visual semantics not fully disentangled and their relationship indirectly modeled. This leads to underutilization of visual information in pre-training, motivating the need for a more balanced approach that explicitly incorporates visual semantics to improve model performance.

To explicitly disentangle the pathology and anatomy semantics and properly model their relationships, we propose a semantic vision-language alignment pipeline: MeDSLIP, **M**edical **D**ual-Stream **L**anguage-**I**mage **P**re-training, and apply it to chest X-ray analysis. MeDSLIP proposes a) a dual-stream mechanism with a disentanglement module to disentangle intertwined pathology and anatomy semantics in images, and b) an interaction modeling module with two contrastive losses to model the relationships between pathology and anatomy semantics. Our contributions are three-fold:

*Firstly*, our dual-stream mechanism separately encodes pathology and anatomy semantics in both medical images and associated reports. In text processing, MeDSLIP extracts pathology and anatomy semantics and prompts them with prior knowledge from humans [23]. In image processing, we disentangle pathology and anatomy semantics from raw images using a disentanglement module. The disentangled visual and textual semantics are then aligned within the pathology-related and anatomy-related streams. By disentangling pathology and anatomy semantics and aligning them in separate streams, MeDSLIP provides a clear understanding of the pathology and anatomy semantics.

*Secondly*, our interaction modeling module exploits a Prototypical Contrastive Loss (ProtoCL) and an Intra-image Contrastive Loss (ICL) to model the relationships between pathology and anatomy semantics. ProtoCL models semantic interactions by aligning the cross-modal, cross-stream information (e.g., pathology in images and anatomy in texts, anatomy in texts and pathology in images). ICL models the pathology-anatomy interactions in images by measuring the co-existence of pathology and anatomy semantics. For example, for a sentence like “Opacity is observed on the bilateral lungs.” in the report and its corresponding image, ProtoCL aligns “opacity” in the image and “lung” in the text, and vice versa, while ICL regularizes the co-existence of “opacity” and “lung” in the image. By modeling these interactions, MeDSLIP captures the rich semantic interactions between pathology and anatomy in both images and texts, leading to a better understanding of the relationships between pathology and anatomy semantics.

*Finally*, to validate MeDSLIP’s effectiveness, we conduct a comprehensive evaluation of classification, grounding, and segmentation tasks under both zero-shot and fine-tuning settings with NIH CXR14 [25], RSNA Pneumonia [26], SIIM-

ACR Pneumothorax [27], and COVIDx CXRv4 [28] datasets. The results demonstrate MeDSLIP’s superior generalizability and transferability. We also conduct an ablation study and qualitative experiments to demonstrate the effectiveness of MeDSLIP and the contributions of its key modules.

## II. METHODOLOGY

Figure 1 shows MeDSLIP’s pipeline. We first disentangle the pathology and anatomy semantics from images and texts and encode them in two distinct streams. Then, the disentangled visual and textual semantics are aligned within each stream. An interaction modeling module is proposed to model the relationship between pathology and anatomy semantics.

### A. Text Processing

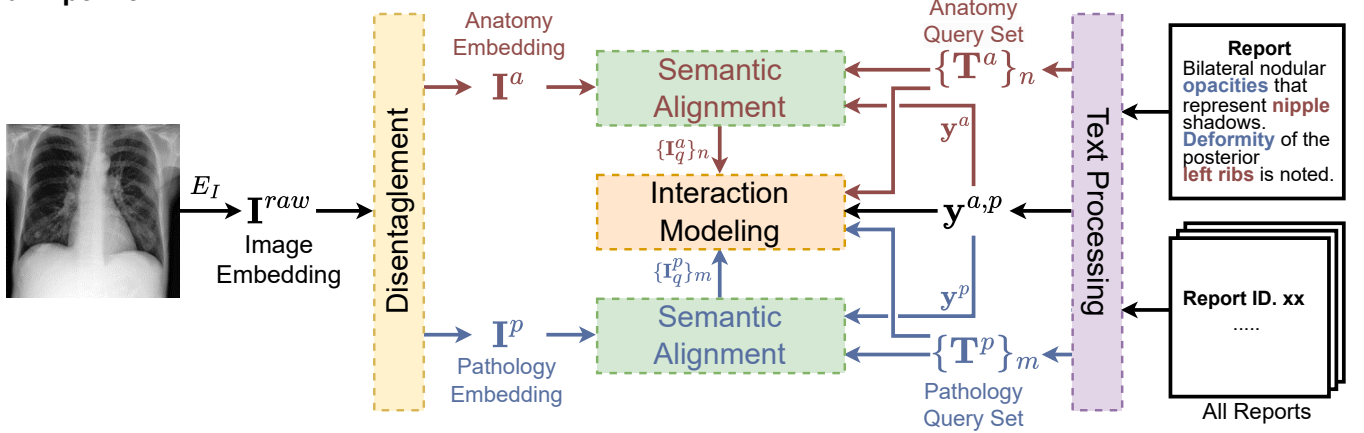
Figure 1b shows the text processing in MeDSLIP. We process all medical reports before pre-training in three steps [23]: (1) Triplet extraction: We extract (pathology, anatomy, existence) triplets from each report [29]. Then, we select the most commonly occurring triplets in the whole dataset to formulate a pathology query set with  $m$  pathology concepts and an anatomy query set with  $n$  anatomy concepts. (2) Prompting: We prompt pathology semantics with domain knowledge derived from professional medical knowledge bases and reliable Internet resources, which provide definitions and explanations of pathology semantics in plain language [23]. We prompt anatomy semantics by a fixed prompting template: “It is located at [ANATOMY]”. (3) Text encoding: We then encode prompted pathology and anatomy query sets by the text encoder, which consists of a frozen pre-trained medical language model [30]  $E_T$  alongside a learnable linear projection layer  $h_T$ . The generated text embeddings are used as queries in semantic alignment and positive/negative samples in interaction modeling, as depicted in Fig. 1d and 1e.

In this context, we obtain three outputs: a pathology query set  $\{\mathbf{T}^p\}_m$ , an anatomy query set  $\{\mathbf{T}^a\}_n$ , and a set of existence label matrices  $\{\mathbf{y}^{a,p}\}$ . Two query sets consist of the text embeddings of the top  $m$  commonly seen pathology semantics and top  $n$  commonly seen anatomy semantics among triplets from all reports in the dataset. The query sets are universal for all reports. The existence matrix  $\mathbf{y}^{a,p}$  is an  $n \times m$  matrix, whose element  $y^{a,p}_{ij}$  indicates whether a pathology observation  $\mathbf{T}^p_j$  exists at a specific anatomy location  $\mathbf{T}^a_i$ . Columns and rows of  $\mathbf{y}^{a,p}$  correspond to pathology and anatomy semantics in query sets. For example, in the sentence “Deformity of the posterior left ribs is noted”, the pathological observation “deformity” exists on the anatomical location “left ribs”. Thus, the existence label  $y^{a,p}_{\text{left ribs}, \text{deformity}}$  is positive. Otherwise, the existence label  $y^{a,p}_{ij}$  will be negative if it is not mentioned or is mentioned as not existing. The meanings of each element in  $\mathbf{y}^{a,p}$  are universal for all data, but the values are unique for each image.

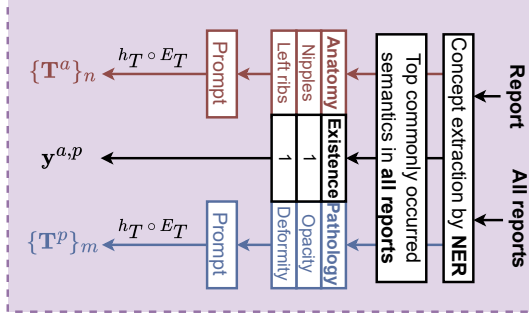
### B. Image Encoding

The image encoding in MeDSLIP’s pre-training comprises three modules: disentanglement, semantic vision-language

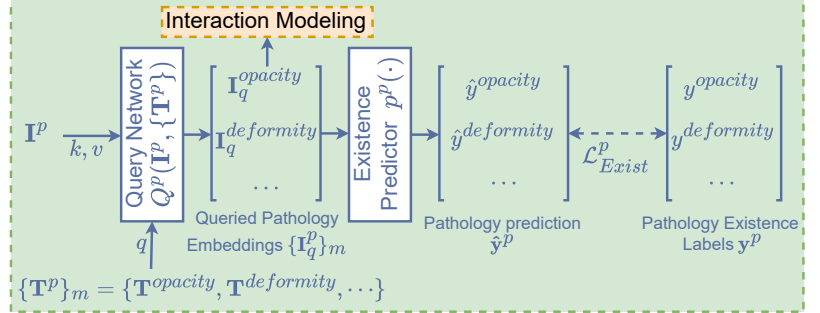
## a. Pipeline



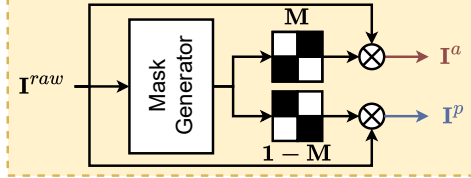
## b. Text Processing Module



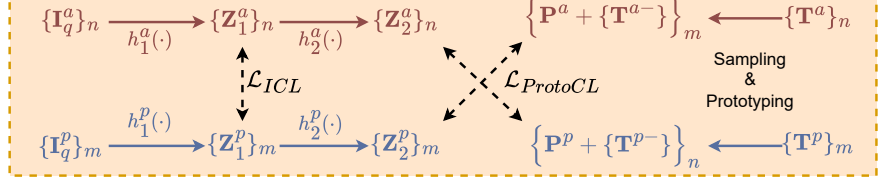
## d. Semantic Alignment in Pathology Stream



## c. Disentanglement Module



## e. Interaction Modeling Module



**Fig. 1: Pipeline of Medical Dual-Stream Language-Image Pre-training (MeDSLIP).** Each module is indicated with a unique color. Symbols with  $\mathbf{I}$  and  $\mathbf{T}$  denote image and text embeddings, respectively.  $Q$  denotes query networks.  $h$  denotes linear projection layers.  $E_I$  and  $E_T$  are image and text encoders, respectively.  $\mathbf{y}$  represents existence labels. The denotations with superscripts  $p$  and  $a$  are pathology-related and anatomy-related. **a. Pipeline:** Reports are processed to extract pathology and anatomy terms, generate text query embedding sets  $\{\mathbf{T}^a\}_n$  and  $\{\mathbf{T}^p\}_m$ , and an existence label matrix,  $\mathbf{y}^{a,p}$ .  $m$  and  $n$  represent that we select top commonly seen  $m$  pathology semantics and  $n$  anatomy semantics in all medical reports. Images are encoded, disentangled, and aligned within corresponding streams. The interaction modeling module regularizes the interactions between pathology and anatomy semantics. **b. Text Processing:** (pathology, anatomy, existence) triplets are extracted from raw reports. Most commonly occurring triplets among all reports are used as query sets, which are prompted and encoded to obtain query embeddings (see Sec. II-A). **c. Disentanglement Module:** It masks raw image embeddings, disentangling pathology and anatomy embedding (Sec. II-B.1). **d. Semantic Alignment:** A query network  $Q^p$  aligns the text query set  $\{\mathbf{T}^p\}_m$  with the image pathology embedding  $\mathbf{I}^p$  and outputs a queried pathology embedding set  $\{\mathbf{I}_q^p\}_m$ . An existence predictor  $p^p$  then checks whether each text semantic exists in the images. A similar alignment process is applied to anatomy semantics (see Sec. II-B.2). **e. Interaction Modeling:**  $\mathcal{L}_{ICL}$  aligns unimodal, cross-stream information, while  $\mathcal{L}_{ProtoCL}$  aligns cross-modal, cross-stream information (see Sec. II-B.3).

alignment, and interaction modeling. A trainable visual encoder  $E_I$  is employed to encode the image into the latent space, where all three modules operate.

**1) Disentanglement Module:** We design a disentanglement module to disentangle the intertwined pathology and anatomy semantics in medical images. As shown in Fig. 1c, we use a mask generator, which takes the raw image embedding  $\mathbf{I}^{raw}$  as input and outputs a mask  $\mathbf{M}$ .  $\mathbf{M}$  has the same shape as the input embedding. The elements in  $\mathbf{M}$  range from 0 to

1. Then, we disentangle pathology and anatomy embeddings  $\mathbf{I}^p$  and  $\mathbf{I}^a$  from raw image embedding  $\mathbf{I}^{raw}$  through element-wise multiplication with  $\mathbf{M}$  and  $1 - \mathbf{M}$ , where  $\mathbf{1}$  is an all-ones matrix with the same size as  $\mathbf{M}$ :

$$\mathbf{I}^p = \mathbf{M} \odot \mathbf{I}^{raw}, \mathbf{I}^a = (\mathbf{1} - \mathbf{M}) \odot \mathbf{I}^{raw}. \quad (1)$$

By disentangling pathology and anatomy semantics into distinct streams, MeDSLIP decouples the intertwined information about characteristics and locations of the diseases, providing

a clearer understanding of the different aspects of diseases.

**2) Semantic Vision-Language Alignment:** After disentanglement, the pathology and anatomy embeddings  $\mathbf{I}^p$  and  $\mathbf{I}^a$  are aligned with corresponding text query embeddings  $\mathbf{T}^p$  and  $\mathbf{T}^a$  via query networks  $Q^p$  and  $Q^a$ , respectively. Since the semantic alignments in the two streams are similar, we use the pathology stream as an example to illustrate. As depicted in Fig. 1d, the query network  $Q^p$  takes two inputs: an image embedding  $\mathbf{I}^p$  and the query set  $\{\mathbf{T}^p\}_m$ . For each query embedding  $\mathbf{T}^p$  in  $\{\mathbf{T}^p\}_m$ ,  $Q^p$  extracts a queried image embedding  $\mathbf{I}_q^p$  from  $\mathbf{I}^p$ :  $\mathbf{I}_q^p = Q^p(\mathbf{I}^p, \mathbf{T}^p)$ . Here,  $\mathbf{I}_q^p$  represents the corresponding features in  $\mathbf{I}^p$  related to the specific query  $\mathbf{T}^p$ . Then,  $\mathbf{I}_q^p$  is passed through a binary existence predictor  $p^p(\cdot)$ , and we obtain a prediction  $\hat{y}^p$  by  $\hat{y}^p = p^p(\mathbf{I}_q^p)$ . The prediction  $\hat{y}^p$  tells whether a specific pathology semantic exists in the image. Repeating the above process for all queries in  $\{\mathbf{T}^p\}_m$ , we obtain a set of pathology predictions  $\hat{\mathbf{y}}^p$ . Finally, we extract the pathology existence labels  $\mathbf{y}^p$  from  $\mathbf{y}^{a,p}$  to calculate the pathology existence loss  $\mathcal{L}_{Exist}^p$ , a binary cross-entropy loss between prediction  $\hat{\mathbf{y}}^p$  and labels  $\mathbf{y}^p$ :

$$\mathcal{L}_{Exist}^p = \mathcal{L}_{BCE}(p^p(Q^p(\mathbf{I}^p, \mathbf{T}^p), \mathbf{y}^p)). \quad (2)$$

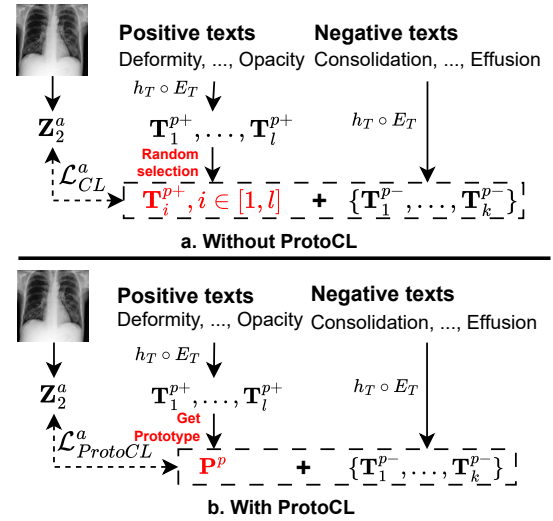
**3) Interaction Modeling:** To properly model the interactions between pathology and anatomy semantics and produce a unified output that retains all relevant information for downstream tasks, we propose an interaction modeling module, as shown in Fig. 1e. This module regularizes the interactions between visual information from one stream and both visual and textual information from the other stream by two specialized losses: prototypical contrastive loss (ProtoCL) and intra-image contrastive loss (ICL). ProtoCL emphasizes the interactions between image embeddings in one stream and text embeddings in the other, while ICL focuses on image embeddings between the two streams.

**a) Prototypical Contrastive Loss (ProtoCL):** ProtoCL regularizes the interactions between cross-modal, cross-stream information by aligning image embeddings in one stream with text embeddings from the other. As is shown in Fig. 2, we use ProtoCL between anatomy image embeddings and pathology text embeddings as an example to illustrate how ProtoCL works and highlight the difference between ProtoCL and conventional contrastive learning.

For medical images and radiology reports, we usually regard existing pathology or anatomy semantics as positive and other unmentioned or non-existing semantics as negative. In this context, there are often multiple positive examples due to the possible coexisting diseases. ProtoCL employs the prototype of all positive examples as the new positive example. The prototype is the center of all positive samples in the textual embedding space, which is calculated by

$$\mathbf{P} = \frac{1}{l} \sum_{i=0}^l \mathbf{T}_i^+. \quad (3)$$

We use a Noise Contrastive Estimation (NCE) loss [31] between prototypes and sampled negatives in ProtoCL, which



**Fig. 2:** Comparison between contrastive learning with or without prototypes, using ProtoCL between anatomy image embeddings and pathology text embeddings as an example. **a.** Conventional contrastive learning without prototypes. **b.** ProtoCL uses the prototype of all positive samples as the new positive example in contrastive learning.

is calculated as follows:

$$\mathcal{L}_{ProtoCL}^a = -\mathbb{E} \left[ \log \frac{\exp(\mathbf{Z}_2^a \cdot \mathbf{P})}{\sum_{i=1}^k \exp(\mathbf{Z}_2^a \cdot \mathbf{T}_i^{p-})} \right]. \quad (4)$$

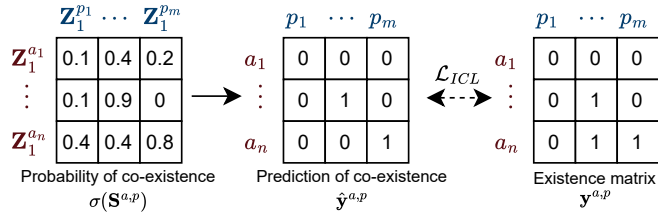
Although conventional contrastive learning can also model the interactions, ProtoCL's key contribution is using prototypes as positive examples. It is motivated by the underutilization of positive information in conventional contrastive learning. When faced with multiple positive examples, conventional contrastive learning tends to randomly select one of the positive examples and leave others unused to keep a low positive-negative ratio [23] for NCE loss, leading to the underutilization of information in the unselected positives. Instead, we design ProtoCL to aggregate all positive information in prototypes without increasing the positive-negative ratio. Therefore, in addition to modeling cross-modal, cross-stream interactions, ProtoCL improves the data efficiency of positive examples by using all positive instances without increasing the positive-to-negative ratio.

**b) Intra-image Contrastive Loss (ICL):** ICL regularizes the visual semantics across two streams by measuring the co-existence of (pathology, anatomy) pairs. A (pathology, anatomy) pair is considered positive if a specific pathology observation is present at a corresponding anatomy location ( $y^{a,p} = 1$ ). Otherwise, the pair is treated as negative ( $y^{a,p} = 0$ ). As depicted in Fig. 3, we first compute a cosine similarity matrix  $\mathbf{S}^{a,p}$  between the queried anatomy image embeddings  $\{\mathbf{I}_q^a\}_n$  and pathology image embeddings  $\{\mathbf{I}_q^p\}_m$ . The element  $s^{a_i,p_j}$  of  $\mathbf{S}^{a,p}$  is given by

$$s^{a_i,p_j} = \langle \mathbf{Z}_1^{a_i}, \mathbf{Z}_1^{p_j} \rangle, i \in [1, n], j \in [1, m], \quad (5)$$

where  $\mathbf{Z}_1^a$  and  $\mathbf{Z}_1^p$  represent linearly projected image embeddings corresponding to specific semantics in the texts. The





**Fig. 3:** Mechanism of intra-image contrastive loss (ICL).  $a_1, \dots, a_n$  and  $p_1, \dots, p_m$  indicating the semantics in query sets. Elements in  $\hat{\mathbf{y}}^{a,p}$  are the predictions, and elements in  $\mathbf{y}^{a,p}$  are ground truths when computing  $\mathcal{L}_{ICL}$ .

similarity matrix  $\mathbf{S}^{a,p}$  is then passed through a sigmoid activation layer  $\sigma(\cdot)$  to produce a probability matrix  $\sigma(\mathbf{S}^{a,p})$ , which represents the likelihood of co-existence for each (pathology, anatomy) pair. Then, we obtain predictions of co-existence  $\hat{\mathbf{y}}^{a,p}$  from the probability matrix  $\sigma(\mathbf{S}^{a,p})$ . We use the existence matrix  $\mathbf{y}^{a,p} = \{l^{a_i, p_j}\}, i \in [1, n], j \in [1, m]$  from the extracted (pathology, anatomy, existence) triplets as ground truths to compute the ICL loss  $\mathcal{L}_{ICL}$ . Each element  $y^{a,p}$  in  $\mathbf{y}^{a,p}$  serves as the existence label for the corresponding (pathology, anatomy) pair. Finally,  $\mathcal{L}_{ICL}$  is obtained by a binary cross-entropy (BCE) loss between  $\hat{\mathbf{y}}^{a,p}$  and  $\mathbf{y}^{a,p}$ :

$$\mathcal{L}_{ICL} = \mathcal{L}_{BCE}(\hat{\mathbf{y}}^{a,p}, \mathbf{y}^{a,p}). \quad (6)$$

The ICL loss encourages the alignment of pathology and anatomy embeddings for positive (pathology, anatomy) pairs while minimizing alignment for negative pairs. This ensures that the model captures fine-grained relationships in images between anatomical structures and pathological conditions, further enhancing interaction modeling from both streams.

In summary, the total loss function in pre-training is

$$\mathcal{L} = \mathcal{L}_{Exist} + \alpha \mathcal{L}_{ProtoCL} + \beta \mathcal{L}_{ICL}, \quad (7)$$

where  $\alpha$  and  $\beta$  are temperature coefficients that balance the scale of the different loss components.  $\mathcal{L}_{Exist}$  and  $\mathcal{L}_{ProtoCL}$  represent the summed losses from the pathology and anatomy streams for existence prediction and prototypical contrastive learning, respectively.

### C. Inference

The existence predictor equips the model with zero-shot classification capability. For instance, when encountering an unseen disease, we prompt the disease in a similar way to how seen pathology semantics are handled during pre-training. The prompted semantics are then encoded by the same text encoder and used as a query in the query network. After obtaining the queried image embedding corresponding to the disease, the existence predictor determines whether the disease is present in the image. Since most downstream tasks are focused on pathology information, we use embeddings in the pathology stream as the model's outputs for downstream tasks ( $\{\mathbf{I}_q^p\}_m$  and  $\hat{\mathbf{y}}^p$  for zero-shot tasks,  $\mathbf{I}^{raw}$  for fine-tuning tasks).

**TABLE I:** Metadata of datasets used in experiments.

Stage	Dataset	# Images	Disease
Pre-training	MIMIC-CXR [2], [32], [33]	377,110	-
	NIH CXR14 [25]	112,000	14 diseases
	RSNA Pneumonia [26]	30,000	Pneumonia
Evaluation	SIIM-ACR Pneumothorax [27]	12,047	Pneumothorax
	COVIDx CXR-4 [28]	84,818	COVID-19

**TABLE II:** Hyperparameters in MeDSLIP pre-training and fine-tuning. The different values in fine-tuning indicate hyperparameters in classification/segmentation fine-tuning tasks.

Name	Pre-training	Fine-tuning
Hidden size	256	256
Image size	224×224	224×224
Batch size	64	64
$n$	50	50
$m$	75	75
$(\alpha, \beta)$	(1, 1)	-
Optimizer	AdamW	AdamW
Epoch	100	100/1000
Peak learning rate	1e-4	1e-4/1e-5
Min learning rate	1e-5	1e-5
Warmup epoch	5	20/50
Warmup learning rate	1e-5	1e-5
Weight decay	0.02	5e-4
Scheduler	cosine annealing	cosine annealing
Time	~ 2 days	minutes to hours
GPU	NVIDIA GeForce RTX4090 24GB	

## III. EXPERIMENT SETTINGS

### A. Datasets

We use MIMIC-CXR [2], [32], [33] for pre-training. NIH CXR14 [25], RSNA Pneumonia [26], SIIM-ACR Pneumothorax [27], and COVIDx CXR-4 [28] are involved in evaluation. The metadata of these datasets is shown in TABLE I.

### B. Implementation

To demonstrate the feasibility of our proposed pipeline, we choose simple structures for each module. MeDSLIP uses ResNet-50 [34] as the image encoder and Bio-ClinicalBERT [30] with a learnable single-layer perceptron (SLP) as the text encoder. For the disentanglement module, an SLP layer is used as the mask generator. We list hyperparameters of MeDSLIP in pre-training and downstream tasks in TABLE II.

### C. Baselines

We compare MeDSLIP to eight leading CNN-based models in the field: ConVIRT [11], MedCLIP [16], GLoRIA [17], BioViL [9], CheXzero [10], UniChest [35], MedKLIP [23], and CXR-CLIP [12]. Since GLoRIA [17] is trained on an in-house dataset, we quote the results in [23]. For BioViL [9], CheXzero [10], UniChest [35], and CXR-CLIP [12], we use the officially released pre-trained weights. MedCLIP [16] doesn't release its code, so we reproduce it according to its paper and pre-train it under our hyperparameters. Besides,

**TABLE III:** Zero-shot classification evaluation with SOTA CNN-based models. For CXR14, metrics refer to the macro average on the 14 diseases.

Models	NIH CXR14 [25]			RSNA Pneumonia [26]			SIIM-ACR Pneumothorax [27]			COVIDx CXR-4 [28]		
	AUROC $\uparrow$	F1 $\uparrow$	ACC $\uparrow$	AUROC $\uparrow$	F1 $\uparrow$	ACC $\uparrow$	AUROC $\uparrow$	F1 $\uparrow$	ACC $\uparrow$	AUROC $\uparrow$	F1 $\uparrow$	ACC $\uparrow$
GLoRIA <sup>+</sup> [17]	66.10	17.32	77.00	71.45	49.01	71.29	53.42	38.23	40.47	-	-	-
MedCLIP [16]	62.91	14.29	54.28	81.66	56.77	73.77	67.94	45.51	61.03	68.84	<u>69.60</u>	64.57
ConVIRT [11]	65.78	15.62	64.83	81.54	56.83	75.83	68.73	45.66	55.53	66.06	66.71	50.04
BioViL [9]	73.92	23.56	79.62	82.54	56.62	76.73	73.53	50.63	65.15	69.82	67.20	54.41
CheXzero [10]	75.65	22.67	<u>86.22</u>	<u>86.11</u>	55.43	66.32	77.28	47.94	56.24	71.56	60.84	69.78
UniChest <sup>†</sup> [35]	-	-	-	79.04	53.55	70.32	89.20	<u>56.52</u>	83.71	68.64	67.47	55.01
CXR-CLIP [12]	75.49	18.17	66.71	82.91	57.57	72.93	84.02	50.10	57.36	71.88	62.57	68.38
CXR-CLIP <sup>†</sup> [12]	-	-	-	83.41	53.80	64.39	85.86	42.28	40.15	54.17	52.38	57.05
MedKLIP [23]	<u>78.00</u>	<u>25.71</u>	85.22	85.27	<u>61.13</u>	<u>79.91</u>	<u>89.60</u>	56.39	<u>84.84</u>	<u>82.77</u>	66.82	<u>73.09</u>
MeDSLIP (Ours)	<b>80.34</b>	<b>29.96</b>	<b>88.93</b>	<b>86.49</b>	<b>63.81</b>	<b>80.98</b>	<b>89.70</b>	<b>58.22</b>	<b>85.05</b>	<b>83.41</b>	<b>77.15</b>	<b>77.21</b>

<sup>+</sup> Because GLoRIA is trained on in-house data, we quote its results in [23]. Because we use a different version of COVIDx [28] than [23], we don't report GLoRIA's results on COVIDx. <sup>†</sup> They include CXR14 as a part of pre-training data. Thus, we don't perform the zero-shot classification task on CXR14.

since MedKLIP [23] doesn't release pre-trained weights, we pre-train it with their official code.

#### D. Metrics

We report the area under the ROC curve (AUROC), F1 score, and accuracy (ACC) for classification tasks, intersection over union (IoU), Dice coefficient, and pointing game score (PG) [23] for grounding tasks, and Dice coefficient for segmentation tasks. Besides, we use *t*-SNE distribution and UMAP to visualize the feature space.

### IV. EXPERIMENT RESULTS

To assess MeDSLIP's capabilities as a vision-language pre-training framework, we evaluate its generalizability and transferability on classification, grounding, and segmentation tasks under both zero-shot and fine-tuning settings. Additionally, we conduct a comprehensive ablation study and validate the effectiveness of each proposed module. Upward arrows ( $\uparrow$ ) of metrics indicate that higher values are better, and downward arrows ( $\downarrow$ ) indicate that lower values are preferred. The best results are highlighted in **bold**, while the second-best results are underlined. We format all results in tables in percentages.

#### A. Generalizability Evaluation

1) *Zero-shot Classification*: We first demonstrate the strong generalizability of MeDSLIP with zero-shot classification tasks. Our zero-shot classification evaluation is conducted in two settings: unseen datasets with seen diseases, and unseen diseases. Unseen datasets with seen diseases indicate that the evaluation datasets are not seen in the model pre-training, but the diseases are seen. Unseen disease indicates that the types of diseases in the evaluation dataset are novel for the model.

a) *Zero-shot Classification on Seen Diseases*: We use NIH CXR14 [25], RSNA Pneumonia [26], and SIIM-ACR Pneumothorax [27] for the zero-shot classification of seen diseases but unseen datasets. In this experiment, although the diseases in the evaluation datasets are included in the pre-training dataset (MIMIC-CXR), the data are novel to MeDSLIP. The results are reported in TABLE III. MeDSLIP demonstrates

strong efficacy and generalizability in zero-shot classification on unseen datasets, outperforming all other baselines across all datasets and metrics by a significant margin.

Because NIH CXR14 [25] contains multiple diseases, including pneumonia and pneumothorax, we use it for visualization and in-depth analysis. NIH CXR14 [25] dataset is a multi-label classification dataset and consists of 14 binary classification problems.

We first present a disease-wise analysis. Figure 4 shows disease-wise AUROC scores for MeDSLIP and seven other baselines. Figure 5 shows disease-wise UMAPs of  $\mathbf{I}_q^p$  between patients of a specific disease (colored points) and healthy controls (gray points). From the disease-wise AUROC bar graph, MeDSLIP outperforms other baselines for most conditions. Particularly, it significantly improves performance for diseases with traditionally low AUROC scores in baselines, such as emphysema (Emp.) and hernia (Her.). This demonstrates MeDSLIP's ability to handle a variety of cardiovascular diseases in chest X-rays. For diseases such as consolidation (Con.), cardiomegaly (Car.), and effusion (Eff.), the UMAPs show less overlap between diseased and healthy embeddings. This corresponds to higher AUROC scores, confirming that better feature separation correlates with better classification performance.

During the evaluation, 14 pathology embeddings  $\{\mathbf{I}_q^p\}_{14}$  are generated by the pathology query network  $Q^p$  for each image, with each embedding  $\mathbf{I}_q^p$  corresponding to one particular disease. Therefore, we visualize the feature space of  $\{\mathbf{I}_q^p\}_{14}$  to further explore how well MeDSLIP is to recognize different diseases. Figure 6 illustrates the label distribution of the dataset and the *t*-SNE distribution of the queried pathology embeddings  $\{\mathbf{I}_q^p\}_{14}$ . The *t*-SNE visualization shows well-clustered embeddings with clear margins between diseases, indicating that MeDSLIP learns a robust feature space even for unseen data distributions during pre-training.

b) *Zero-shot Classification on Unseen Diseases*: We evaluate MeDSLIP's zero-shot classification ability on unseen diseases using the COVIDx CXR-4 [28] dataset. The pre-training dataset, MIMIC-CXR [2], collected data between 2011 and 2016. It predated the COVID-19 pandemic, which began

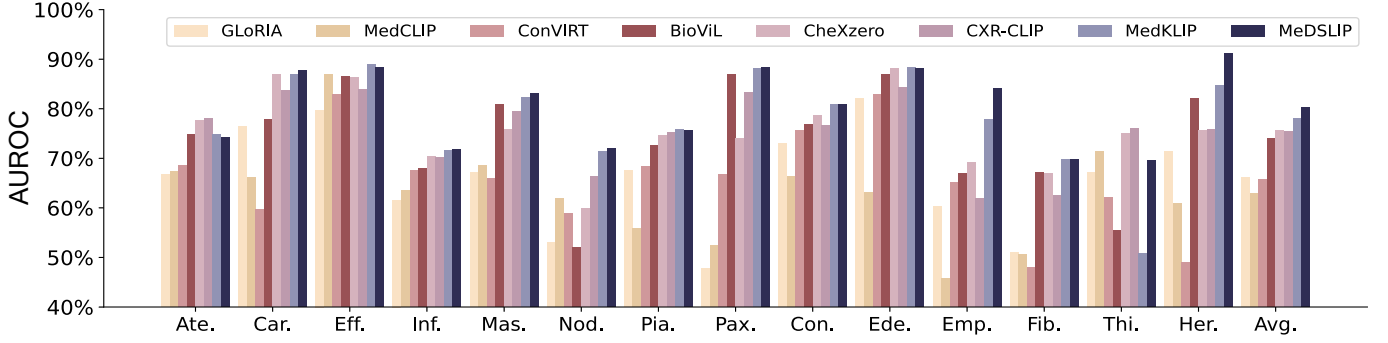


Fig. 4: Disease-wise AUROCs of zero-shot classification on NIH CXR14 dataset [25] show MeDSLIP outperforms other baselines on most of the diseases. AUROCs are calculated between the positive patients of each disease and other health controls across all data.

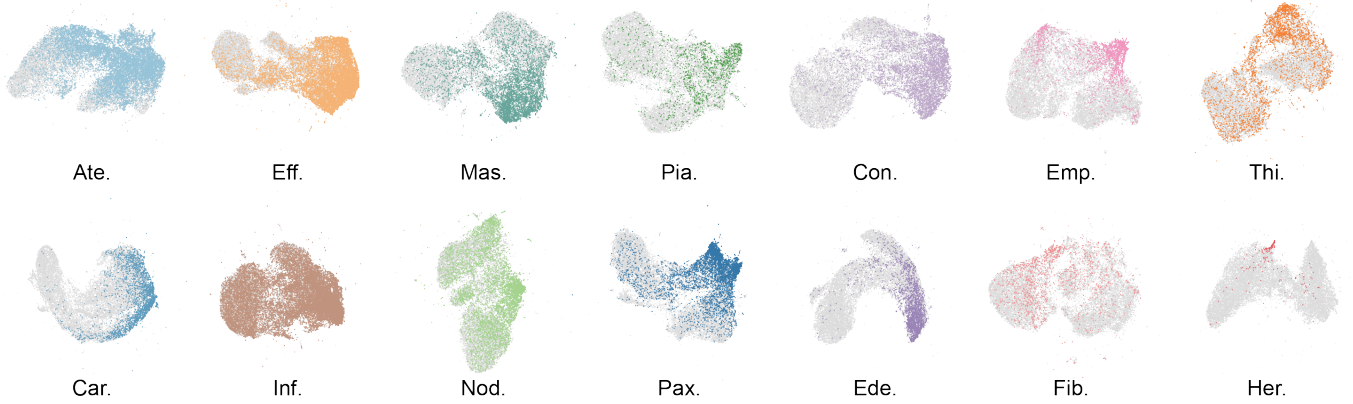


Fig. 5: Disease-wise UMAPs of  $\{I_q^p\}_{14}$  of NIH CXR14 dataset [25]. Gray points in each UMAP represent  $I_q^p$  of healthy controls, while colored points denote  $I_q^p$  of patients with the corresponding disease. The diseases with higher AUROC scores in TABLE III tend to be more distinct and well-clustered.

in 2019. Therefore, COVID-19 is a novel disease class for evaluating MeDSLIP.

Since COVID-19 is not included in previous prompts, we designed a prompt: “COVID-19, caused by the SARS-CoV-2 virus, primarily affects the respiratory system and can be identified on chest X-rays by characteristic bilateral, peripheral ground-glass opacities. These radiographic findings are most commonly seen in the lower lobes of the lungs and can progress to multifocal consolidation as the disease advances.”

We present the results of zero-shot classification on the COVIDx CXR-4 dataset [28] in TABLE III. In the zero-shot classification task, MeDSLIP consistently outperforms all state-of-the-art (SOTA) methods across all metrics. Notably, MeDSLIP improves accuracy by at least 4.12% and increases the F1 score by 7.55%, demonstrating its robustness and superior generalizability in handling unseen diseases.

2) *Zero-shot Grounding*: MeDSLIP’s grounding ability is evaluated through a zero-shot grounding task on the RSNA Pneumonia dataset [26]. We use attention maps from the pathology query network  $Q^p$  and a predefined threshold to identify the abnormal regions. Regions with attention values exceeding the threshold are regarded as abnormal regions, which are then compared against the ground truth to compute the evaluation metrics.

Four models are included in this experiment, with results

TABLE IV: Zero-shot grounding tasks on RSNA Pneumonia dataset [26].

Models	Dice $\uparrow$	IoU $\uparrow$	PG $\uparrow$
GLoRIA [17]	34.68	21.82	76.07
BioViL [9]	44.34	30.76	84.12
MedKLIP [23]	49.63	34.32	87.02
MeDSLIP	<b>50.60</b>	<b>35.47</b>	<b>91.10</b>

presented in TABLE IV. MeDSLIP outperforms all baselines across the three metrics, demonstrating its strong grounding performance. It achieves the highest Dice coefficient of 50.60%, indicating superior overlap between predicted regions and ground truths. The model also records the highest IoU of 35.47%, demonstrating better performance in grounding tasks. Especially, MeDSLIP achieves a Pointing Game Score of 91.10%, significantly outperforming the other models. This result highlights MeDSLIP’s superior capability to localize objects within the images. Together, these findings confirm MeDSLIP’s overall superiority in zero-shot grounding tasks compared to competing methods.

We visualize attention maps alongside annotated bounding boxes in Fig. 7. In pneumonia grounding tasks, MeDSLIP first localizes the lungs, assigning them more attention than



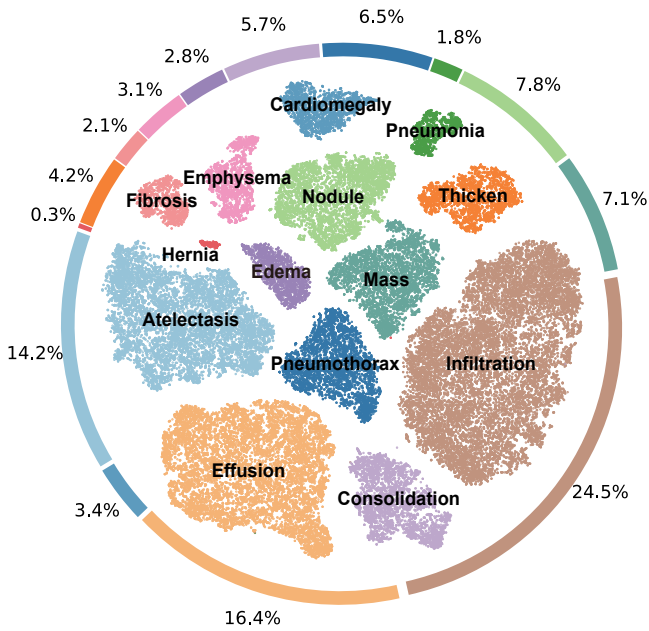


Fig. 6: MeDSLIP clearly distinguishes different pathology semantics. The center of the figure is the  $t$ -SNE distribution of queried pathology embeddings  $\{\mathbf{I}_q^p\}_m$  to 14 diseases. The outside donut chart shows the class distribution in the NIH CXR14 dataset [25].

TABLE V: Fine-tuning on classification and segmentation tasks.

Task	Dataset	Model		
		Data Ratio	MedKLIP [23]	MeDSLIP
Classification (AUROC $\uparrow$ )	CXR14 [25]	1%	66.20	<b>70.60</b>
		10%	76.66	<b>77.20</b>
		100%	80.30	<b>80.50</b>
	SIIM-ACR [27]	1%	86.27	<b>88.06</b>
		10%	89.90	<b>90.96</b>
		100%	93.10	<b>93.83</b>
	COVIDx [28]	1%	65.02	<b>67.41</b>
		10%	74.47	<b>74.64</b>
		100%	77.46	<b>80.38</b>
Segmentation (Dice $\uparrow$ )	SIIM-ACR [27]	1%	67.71	<b>70.39</b>
		10%	75.52	<b>77.53</b>
		100%	76.96	<b>77.75</b>

surrounding areas. It then highlights abnormal regions within the lungs with significantly stronger attention responses.

## B. Transferability Evaluation

1) *Fine-tuning Classification*: We evaluate MeDSLIP's transferability on fine-tuning classification tasks using the NIH CXR14 [25], SIIM-ACR Pneumothorax [27], and COVIDx CXR-4 [28] datasets. Classification models typically comprise a feature extractor to generate embeddings and a classification head to predict. We employ the pre-trained image encoder of MeDSLIP as the feature extractor and a randomly initialized binary classifier as the classification head. During fine-tuning, the encoder is frozen while the classifier is trainable. The model is fine-tuned using 1%, 10%, and 100% of the training

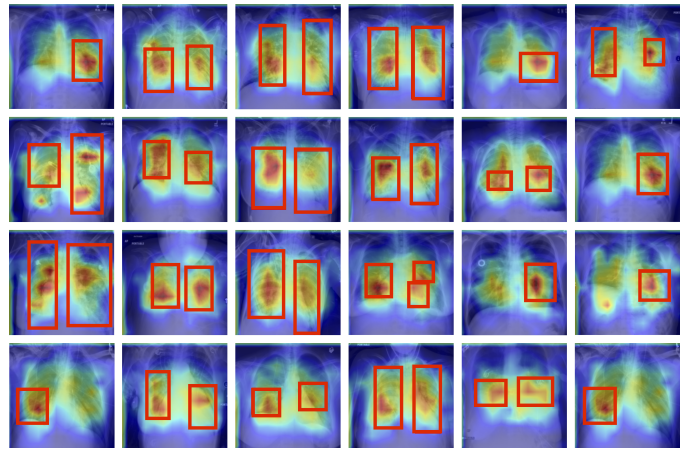


Fig. 7: Visualization of attention maps and ground truths of zero-shot grounding task on RSNA Pneumonia dataset [26] shows that MeDSLIP identifies abnormal regions. The red boxes are ground truths.

set. The AUROC scores are reported in TABLE V. Across different ratios of training data, MeDSLIP consistently achieves higher AUROC scores than the baseline, demonstrating its superiority in fine-tuning classification tasks.

Notably, MeDSLIP shows a significant advantage when trained on the smallest data size, highlighting its data efficiency with limited data. While the performance gap narrows as the training data size increases, MeDSLIP still maintains a clear lead, underscoring its robustness and adaptability in classification tasks, even in data-rich scenarios.

2) *Fine-tuning Segmentation*: We further evaluate MeDSLIP on fine-tuning segmentation tasks using the SIIM-ACR Pneumothorax dataset [27]. Segmentation models typically comprise an encoder to extract features from raw images and a pixel-dense decoder to generate segmentation maps. For this experiment, we use ResUNet [37] as the backbone network. The pre-trained image encoder from MeDSLIP is employed to initialize the encoder of ResUNet, while the decoder is randomly initialized. During fine-tuning, the encoder remains frozen, and only the decoder is trainable. The training data sizes used for segmentation experiments mirror those of the fine-tuning classification tasks, with 1%, 10%, and 100% of the SIIM-ACR Pneumothorax dataset [27]. The dice scores for these experiments are presented in TABLE V.

The results reveal a consistent trend similar to that observed in the classification tasks. Across all training data sizes, MeDSLIP outperforms the baseline. Notably, the performance improvement of MeDSLIP is most pronounced when less training data is available, demonstrating its high data efficiency with limited data. These findings further illustrate the robustness and superiority of MeDSLIP in fine-tuning segmentation tasks, particularly in data-scarce scenarios.

## C. Ablation Study

To show the contributions of each module in MeDSLIP, we perform an ablation study, with results presented in TABLE VI. This study evaluates the roles of ProtoCL and ICL within the interaction modeling module separately.



**TABLE VI:** Ablation study under zero-shot setting: BL, PCL, ICL, DIS, represent baseline (MedKLIP), ProtoCL loss, ICL loss, and disentanglement module with dual-stream structure and mask generator.

Exp. ID	Modules				Classification						Grounding <sup>♦</sup>	
					CXR14 <sup>†</sup> [25]			RSNA [26]			RSNA [26]	
	BL	DIS	PCL	ICL	AUROC <sup>↑</sup>	F1 <sup>↑</sup>	ACC <sup>↑</sup>	AUROC <sup>↑</sup>	F1 <sup>↑</sup>	ACC <sup>↑</sup>	Point	Score <sup>↑</sup>
1	✓				76.15	22.98	84.08	85.74	61.88	79.76	87.02	
2	✓	✓			77.30	25.28	85.00	85.62	62.13	78.60	82.45	
3	✓		✓		77.24	23.42	86.68	<u>85.77</u>	<u>62.57</u>	<u>81.51</u>	<u>87.69</u>	
4	✓	✓	✓		<u>77.81</u>	<b>26.59</b>	86.90	84.99	60.64	80.08	83.98	
5	✓	✓	✓	✓	<b>78.14</b>	<u>26.01</u>	<b>88.88</b>	<b>86.49</b>	<b>63.18</b>	<b>82.97</b>	<b>91.10</b>	

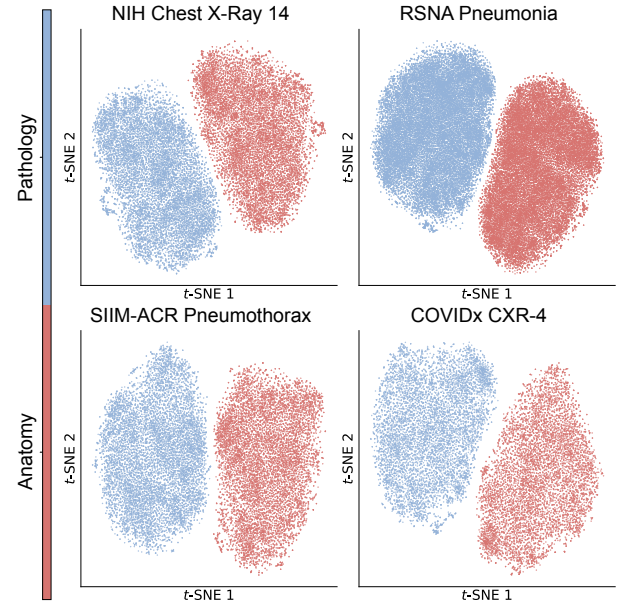
♦ We use the hyperparameter-agnostic pointing game score [36] in the grounding task to avoid the effects of hyperparameter selection. † The study on CXR14 [25] uses the full dataset.

**1) Effect of Disentanglement Module:** The disentanglement module clearly disentangles pathology and anatomy semantics. Comparisons between Exp. 1, 2, and Exp. 3, 4 in TABLE VI show a drop in grounding performance while maintaining competitive or even improved classification performance. This observation aligns with theoretical expectations: Without comprehensive information exchange in interaction modeling, the disentanglement module separates pathology and anatomy semantics, and we should only use the pathology information in downstream tasks. Consequently, model experiences reduced grounding performance unless the separated anatomy information is re-integrated via the interaction modeling module.

As mentioned in Sec. II-B.1, pathology embeddings capture the types of diseases that are essential for classification tasks, while anatomy embeddings describe the locations of diseases, which are important for grounding tasks. We separate information about two aspects of diseases into two distinct streams by disentangling pathology and anatomy semantics, and only use the pathology outputs for downstream tasks. In this context, the pathology information remains intact or becomes more distinct, resulting in competitive classification performance. However, with only a part of or even no cross-stream interaction modeling, the anatomy (location) information is only partially included or not included in the outputs, leading to a decrease in grounding performance.

To further explore the quality of disentanglement, we visualize the disentangled pathology and anatomy embeddings  $I^a$  and  $I^p$  in  $t$ -SNE space, as shown in Fig. 8. The visualization shows clear clusters for pathology and anatomy information across evaluation datasets, confirming that the module successfully disentangles these semantics into their respective streams.

**2) Effect of ProtoCL:** ProtoCL models cross-modal, cross-stream interaction and improve the data efficiency of positive samples. Comparing Exp. 1 with 3 in TABLE VI, the model with ProtoCL achieves superior performance compared to conventional contrastive learning. This demonstrates improved data efficiency for positive samples when information is not disentangled. In the case of disentangled information (Exp. 2 and 4), ProtoCL improves grounding performance by developing the relationship between anatomy information with pathology outputs. These results confirm that ProtoCL fulfills its design objectives of improving positive data efficiency and



**Fig. 8:**  $t$ -SNE graphs of  $I^a$  and  $I^p$  in evaluation datasets indicate the disentanglement module separates pathology and anatomy semantics. The red points are anatomy embeddings, and the blue ones are pathology embeddings.

modeling cross-stream interaction properly.

**3) Effect of ICL:** ICL further facilitates cross-stream interaction modeling by regularizing disentangled visual pathology and anatomy semantics. Comparing Exp. 4 and 5 in TABLE VI, MeDSLIP with ICL demonstrates clear improvements across almost all metrics and datasets. Notably, the model's grounding performance with both ICL and disentanglement (Exp. 5) outperforms that of models without these components (Exp. 1 and 3). These results indicate that ICL successfully regularizes the interaction modeling between the disentangled pathology and anatomy information in a more organized manner, further enhancing the model's ability to perform grounding and classification tasks effectively.

## V. CONCLUSION

To address the entanglement of pathology and anatomy semantics and properly model their relationships, we propose

a semantic vision-language alignment pipeline: MeDSLIP, **Medical Dual-Stream Language-Image Pre-training**. It explicitly disentangles pathology and anatomy semantics in texts and images to enhance the model's ability to utilize these semantics effectively and properly model their interactions. MeDSLIP demonstrates superior generalizability and transferability by consistently outperforming eight other baselines in all experiments. With its superiority, MeDSLIP offers a robust tool to aid in complex clinical tasks and lays the groundwork for future innovations in AI-driven healthcare.

## REFERENCES

- [1] E. Samei and D. Peck, *Anatomy, Physiology, and Pathology in Imaging*, 02 2019, pp. 55–87.
- [2] A. E. Johnson, T. J. Pollard, S. J. Berkowitz, N. R. Greenbaum, M. P. Lungren, C.-y. Deng, R. G. Mark, and S. Horng, “MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports,” *Scientific Data*, vol. 6, no. 1, pp. 317–324, 2019.
- [3] W. Chen, P. Wang, H. Ren, L. Sun, Q. Li, Y. Yuan, and X. Li, “Medical image synthesis via fine-grained image-text alignment and anatomy-pathology prompting,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2024, pp. 240–250.
- [4] A. Jaus, C. Seibold, S. Reiß, L. Heine, A. Schily, M. Kim, F. H. Bahnsen, K. Herrmann, R. Stiefelhagen, and J. Kleesiek, “Anatomy-guided pathology segmentation,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2024, pp. 3–13.
- [5] I. Hartsock and G. Rasool, “Vision-language models for medical report generation and visual question answering: A review,” *arXiv preprint arXiv:2403.02469*, 2024.
- [6] C. Liu, Y. Jin, Z. Guan, T. Li, Y. Qin, B. Qian, Z. Jiang, Y. Wu, X. Wang, Y. F. Zheng *et al.*, “Visual-language foundation models in medicine,” *The Visual Computer*, pp. 1–20, 2024.
- [7] B. Azad, R. Azad, S. Eskandari, A. Bozorgpour, A. Kazerouni, I. Rekik, and D. Merhof, “Foundational models in medical imaging: A comprehensive survey and future vision,” *arXiv preprint arXiv:2310.18689*, 2023.
- [8] C. Chen, C. Qin, H. Qiu, G. Tarroni, J. Duan, W. Bai, and D. Rueckert, “Deep learning for cardiac image segmentation: a review,” *Frontiers in Cardiovascular Medicine*, vol. 7, 2020.
- [9] B. Boecking, N. Usuyama, S. Bannur, D. C. Castro, A. Schwaighofer, S. Hyland, M. Wetscherek, T. Naumann, A. Nori, J. Alvarez-Valle *et al.*, “Making the most of text semantics to improve biomedical vision-language processing,” in *European Conference on Computer Vision*. Springer, 2022, pp. 1–21.
- [10] E. Tiu, E. Talus, P. Patel, C. P. Langlotz, A. Y. Ng, and P. Rajpurkar, “Expert-level detection of pathologies from unannotated chest X-ray images via self-supervised learning,” *Nature Biomedical Engineering*, vol. 6, no. 12, pp. 1399–1406, 2022.
- [11] Y. Zhang, H. Jiang, Y. Miura, C. D. Manning, and C. P. Langlotz, “Contrastive learning of medical visual representations from paired images and text,” in *Machine Learning for Healthcare Conference*. PMLR, 2022, pp. 2–25.
- [12] K. You, J. Gu, J. Ham, B. Park, J. Kim, E. K. Hong, W. Baek, and B. Roh, “CXR-CLIP: Toward large scale chest X-ray language-image pre-training,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2023, pp. 101–111.
- [13] C. Liu, S. Cheng, C. Chen, M. Qiao, W. Zhang, A. Shah, W. Bai, and R. Arcucci, “M-FLAG: Medical vision-language pre-training with frozen language models and latent space geometry optimization,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2023, pp. 637–647.
- [14] H.-Y. Zhou, X. Chen, Y. Zhang, R. Luo, L. Wang, and Y. Yu, “Generalized radiograph representation learning via cross-supervision between images and free-text radiology reports,” *Nature Machine Intelligence*, vol. 4, no. 1, pp. 32–40, 2022.
- [15] Z. Wan, C. Liu, M. Zhang, J. Fu, B. Wang, S. Cheng, L. Ma, C. Quilodrán-Casas, and R. Arcucci, “Med-UniC: Unifying cross-lingual medical vision-language pre-training by diminishing bias,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [16] Z. Wang, Z. Wu, D. Agarwal, and J. Sun, “MedCLIP: Contrastive learning from unpaired medical images and text,” *arXiv preprint arXiv:2210.10163*, 2022.
- [17] S.-C. Huang, L. Shen, M. P. Lungren, and S. Yeung, “GLoRIA: A multimodal global-local representation learning framework for label-efficient medical image recognition,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 3942–3951.
- [18] C. Liu, S. Cheng, M. Shi, A. Shah, W. Bai, and R. Arcucci, “IMITATE: Clinical prior guided hierarchical vision-language pre-training,” *IEEE Transactions on Medical Imaging*, pp. 519–529, 2024.
- [19] F. Wang, Y. Zhou, S. Wang, V. Vardhanabuthi, and L. Yu, “Multi-granularity cross-modal alignment for generalized medical visual representation learning,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 33 536–33 549, 2022.
- [20] X. Zhao, Z.-Y. Liu, F. Liu, G. Li, Y. Dou, and S. Peng, “Report-concept textual-prompt learning for enhancing X-ray diagnosis,” in *ACM Multimedia 2024*, 2024.
- [21] C. Wu, X. Zhang, Y. Wang, Y. Zhang, and W. Xie, “K-Diag: Knowledge-enhanced disease diagnosis in radiographic imaging,” *arXiv preprint arXiv:2302.11557*, 2023.
- [22] X. Zhang, C. Wu, Y. Zhang, W. Xie, and Y. Wang, “Knowledge-enhanced visual-language pre-training on chest radiology images,” *Nature Communications*, vol. 14, no. 1, pp. 4542–4553, 2023.
- [23] C. Wu, X. Zhang, Y. Zhang, Y. Wang, and W. Xie, “MedKLIP: Medical knowledge enhanced language-image pre-training,” *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023.
- [24] V. M. H. Phan, Y. Xie, Y. Qi, L. Liu, L. Liu, B. Zhang, Z. Liao, Q. Wu, M.-S. To, and J. W. Verjans, “Decomposing disease descriptions for enhanced pathology detection: A multi-aspect vision-language pre-training framework,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2024, pp. 11 492–11 501.
- [25] X. Wang, Y. Peng, L. Lu, Z. Lu, M. Bagheri, and R. M. Summers, “ChestX-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2097–2106.
- [26] G. Shih, C. C. Wu, S. S. Halabi, M. D. Kohli, L. M. Prevedello, T. S. Cook, A. Sharma, J. K. Amorosa, V. Arteaga, M. Galperin-Aizenberg *et al.*, “Augmenting the national institutes of health chest radiograph dataset with expert annotations of possible pneumonia,” *Radiology: Artificial Intelligence*, vol. 1, no. 1, 2019.
- [27] A. Zawacki, C. Wu, G. Shih, J. Elliott, M. Fomitchev, M. Hussain, ParasLakhani, P. Culliton, and S. Bao, “Siim-acr pneumothorax segmentation,” 2019. [Online]. Available: <https://kaggle.com/competitions/siim-acr-pneumothorax-segmentation>
- [28] L. Wang, Z. Q. Lin, and A. Wong, “COVID-Net: a tailored deep convolutional neural network design for detection of COVID-19 cases from chest X-ray images,” *Scientific Reports*, vol. 10, no. 1, pp. 19 549–19 560, Nov 2020.
- [29] S. Jain, A. Agrawal, A. Saporta, S. Q. Truong, D. N. Duong, T. Bui, P. Chambon, Y. Zhang, M. P. Lungren, A. Y. Ng *et al.*, “Radgraph: Extracting clinical entities and relations from radiology reports,” *arXiv preprint arXiv:2106.14463*, 2021.
- [30] E. Alsentzer, J. R. Murphy, W. Boag, W.-H. Weng, D. Jin, T. Naumann, and M. McDermott, “Publicly available clinical BERT embeddings,” *arXiv preprint arXiv:1904.03323*, 2019.
- [31] M. U. Gutmann and A. Hyvärinen, “Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics,” *Journal of Machine Learning Research*, vol. 13, no. 11, pp. 307–361, 2012.
- [32] A. E. Johnson, T. J. Pollard, R. G. Mark, S. J. Berkowitz, and S. Horng, “MIMIC-CXR Database (version 2.0.0),” 2019. [Online]. Available: <https://doi.org/10.13026/C2JT1Q>
- [33] A. E. Johnson, T. J. Pollard, N. R. Greenbaum, M. P. Lungren, C.-y. Deng, Y. Peng, Z. Lu, R. G. Mark, S. J. Berkowitz, and S. Horng, “MIMIC-CXR-JPG, a large publicly available database of labeled chest radiographs,” *arXiv preprint arXiv:1901.07042*, 2019.
- [34] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [35] T. Dai, R. Zhang, F. Hong, J. Yao, Y. Zhang, and Y. Wang, “UniChest: Conquer-and-divide pre-training for multi-source chest X-ray classification,” *IEEE Transactions on Medical Imaging*, vol. 43, no. 8, pp. 2901–2912, 2024.

- [36] J. Zhang, S. A. Bargal, Z. Lin, J. Brandt, X. Shen, and S. Sclaroff, "Top-down neural attention by excitation backprop," *International Journal of Computer Vision*, vol. 126, no. 10, pp. 1084–1102, 2018.
- [37] F. I. Diakogiannis, F. Waldner, P. Caccetta, and C. Wu, "ResUNet-a: A deep learning framework for semantic segmentation of remotely sensed data," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 162, pp. 94–114, 2020.