



# Multi-Modal Understanding and Generation for Medical Images and Text via Vision-Language Pre-Training

Jong Hak Moon , Hyungyung Lee , Woncheol Shin, Young-Hak Kim , and Edward Choi 

**Abstract**—Recently a number of studies demonstrated impressive performance on diverse vision-language multi-modal tasks such as image captioning and visual question answering by extending the BERT architecture with multi-modal pre-training objectives. In this work we explore a broad set of multi-modal representation learning tasks in the medical domain, specifically using radiology images and the unstructured report. We propose Medical Vision Language Learner (MedViLL), which adopts a BERT-based architecture combined with a novel multi-modal attention masking scheme to maximize generalization performance for both vision-language understanding tasks (diagnosis classification, medical image-report retrieval, medical visual question answering) and vision-language generation task (radiology report generation). By statistically and rigorously evaluating the proposed model on four downstream tasks with three radiographic image-report datasets (MIMIC-CXR, Open-I, and VQA-RAD), we empirically demonstrate the superior downstream task performance of MedViLL against various baselines, including task-specific architectures.

**Index Terms**—Healthcare, medical, multimodal learning, representation learning, self-supervised learning, vision-and-language.

## I. INTRODUCTION

**V**ISION-LANGUAGE (VL) multi-modal research using radiographic images and associated free-text description (e.g., Chest X-rays and radiology report) is one of the

most important and interesting works in the medical informatics [4], [6], [18], [21], [22], [29], [37]. Although each VL modality provides different representations to the researcher, images and reports contain mutually helpful semantic information. Consequently, advances in VL multi-modal research can be beneficial in improving the quality of clinical care by providing automated support for a variety of tasks such as diagnosis classification [6], [18], [37], report generation [22], [37]. Owing to the high dimensionality, heterogeneity, and systemic biases, however, handling both image and clinical report to learn joint representation poses significant technical challenges.

The development of VL multi-modal learning has produced tremendous progress recently by extending the BERT-based architecture [10] in deep learning area. BERT-based VL model is the typical pretrain-then-transfer approach that makes the model learn a representation of each modality by performing multiple pre-training tasks. After pre-training the model, it is transferred to various vision language understanding (VLU) (e.g., visual question answering, text-conditioned image retrieval and vice versa) and vision language generation (VLG) (e.g., image captioning) downstream tasks by making only minor additions to the base architecture. However, despite significant improvements reported for a wide range of downstream tasks utilizing pre-trained models, most previous studies focused on either the VLU tasks or the VLG tasks [14], [24], [34], suggesting the challenging nature of learning meaningful representations for both VLU and VLG at the same time. Some recent studies tried to tackle both tasks at the same time by proposing a hybrid model using the encoder and decoder of the Transformer [32], [33], [39], or with a unified BERT model by sharing knowledge using different types of the self-attention mask [41]. These works demonstrated promising results even when a single BERT-based architecture was trained to aim at both tasks.

While VL multi-modal pre-training has no doubt seen significant progress in recent years, it was mainly developed under the context of general domain (e.g., using MS-COCO). Vision and language, however, is one of the most frequently used information in the medical domain as well, often produced in the form of radiology images and corresponding free-text report. VL multi-modal pre-training therefore has great potential to be widely used in healthcare such improving diagnosis accuracy, automatically generating reports, or answering questions from physicians. Despite its huge potential, VL multi-modal pre-training in the

Manuscript received 3 December 2021; revised 19 May 2022 and 8 August 2022; accepted 7 September 2022. Date of publication 19 September 2022; date of current version 6 December 2022. This work was supported in part by Samsung Electronics under Grant IO201211-08109-01, in part by the Institute of Information and Communications Technology Planning and Evaluation (IITP) under Grant 2019-0-00075, in part by Artificial Intelligence Graduate School Program (KAIST), and in part by the National Research Foundation of Korea (NRF) under Grant NRF-2020H1D3A2A03100945 funded by the Korea government (MSIT). (Jong Hak Moon and Hyungyung Lee contributed equally to this work.) (Corresponding author: Edward Choi.)

Jong Hak Moon, Hyungyung Lee, Woncheol Shin, and Edward Choi are with the Graduate School of AI, KAIST, Daejeon 34141, South Korea (e-mail: jhak.moon@kaist.ac.kr; ttumyche@kaist.ac.kr; swc1905@kaist.ac.kr; edwardchoi@kaist.ac.kr).

Young-Hak Kim is with the Department of Cardiology, Asan Medical Center, University of Ulsan College of Medicine, Seoul 05505, South Korea (e-mail: mdyhkim@amc.seoul.kr).

The source code is publicly available at: <https://github.com/SuperSupermoon/MedViLL>.

Digital Object Identifier 10.1109/JBHI.2022.3207502

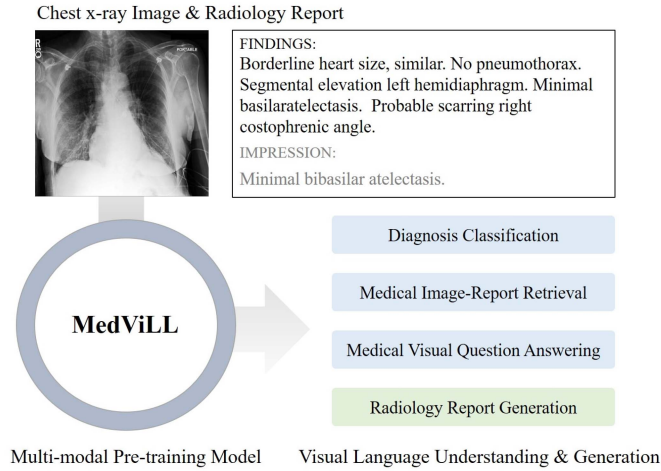


Fig. 1. Overview of MedViLL. During the pre-training, MedViLL learns joint representation, then fine-tuned for VLU and VLG tasks.

medical domain has only recently received attention, where Li et al. [20] only demonstrated improved diagnosis accuracy of VL pre-trained models. In order to truly understand, however, whether a model has effectively learned both vision and text representation, it must be evaluated on diverse VLU and VLG tasks beyond simple diagnosis classification. This motivates us to investigate whether an integrated model is possible for a wide range of VLU and VLG tasks.

In this paper, we aim to develop a model that can learn multipurpose joint representations of vision and language in the medical domain (Fig. 1). The more immediate focus of our approach to enhance downstream tasks such as diagnosis and treatment delivery is case-based reasoning, discovering underlying patterns in data, and generating semantically accurate disease profiles. The main contributions of this paper can be summarized as follows:

- 1) We propose Medical Vision Language Learner (MedViLL), a multi-modal pre-training model for medical images and reports with a novel self-attention scheme.
- 2) We demonstrate the effectiveness of our approach with detailed ablation study on extensive vision-language understanding and generation-based downstream tasks, including diagnosis classification, medical image-report retrieval, medical visual question answering, and radiology report generation.
- 3) We demonstrate the generalization ability of our approach under the transfer learning setting using two separate Chest X-ray datasets, where we pre-train a model on one dataset and perform diverse downstream tasks on another.

To the best of our knowledge, this is the first study that conducts both VLU and VLG tasks with a unified VL pre-training model in the medical domain. We expect that our pretrained VL model will enable more effective cross-task knowledge sharing, and reduce the development costs by eliminating the need for separate models for different tasks.

## II. RELATED WORK

### A. Radiology Practices

In radiology practice, physicians identify various clinical findings based on radiographic images and the patient's clinical history, then summarize these findings and overall impressions in a clinical report [17], [30]. Diagnostic observations are described as positive, negative, or uncertain about the clinical findings, including the detailed location and severity of the findings. Such clinical reports are currently being used as a standard communication method in the clinical setting. A combination of vision and language data helps further improve the model performance in both image annotation and automatic report generation [21].

### B. VL Multimodal Researches in the Medical Domain

Although various models have been gradually developed for language modeling [2], [19], [23], [31], CNN-RNN based models still dominate in VL multi-modal learning in the medical domain, and these models were mainly designed for a task-specific method of either VLU or VLG tasks. TieNet [37] is a pioneering CNN-RNN model with image-report attention mechanism for VLU (e.g., diagnosis classification) and VLG (e.g., report generation) tasks by using ChestX-ray14 [36] dataset. Liu et al. [22] only focus on the VLG task to generate the radiology report utilizing a CNN-RNN-RNN architecture with a hierarchical generation strategy from the MIMIC-CXR [16] and Open-I dataset [9]. Hsu et al. [12] focus on a VLU task, specifically image-report retrieval in the MIMIC-CXR dataset, based on supervised and unsupervised methods. The most recent studies [20], [4], [40] focus on either VLU or VLG task. Li et al. [20] compares 4 different BERT-based pre-training models on a VLU task, specifically classifying thoracic findings in the MIMIC-CXR and Open-I dataset. With a focus on VLG task [4], [40], EMIXER [4] is a GAN-based approach that simultaneously generates a pair of X-ray images and corresponding reports based on diagnosis labels. Yang et al. [40] proposed MedWriter focusing on report generation task that incorporates a hierarchical retrieval mechanism to automatically extract both report and sentence-level templates. In this paper, we focus on learning a joint representation of a single image and its corresponding report to perform both VLU and VLG tasks with fine-tuning.

### C. VL Multimodal Researches in the General Domain

For better understanding of VL multimodality, many works have been proposed recently [7], [14], [24], [33], [34], [41] in the general domain. Among numerous variants of VL pre-training setup, we focus on three components that are most relevant to our approach: input embedding stream, visual feature embedding, and downstream tasks.

**1) Input Embedding Stream:** Existing models can be divided into two groups based on their architecture as a single- [7], [14], [33] or two-stream [24], [34] with the marginal difference in downstream task performance [5], [26]. However,

the two-stream architecture has a greater number of parameters, whereas the single stream architecture allows early interaction between two modalities by sharing parameters and processing stacks [7], [14], [41]. For architectural simplicity and time/space efficiency, we design our model with a single-stream architecture.

**2) Visual Feature Embedding:** For visual feature embedding, most of the recent works [7], [24] are inspired by [3] utilizing pre-trained object detectors [28] to extract the region-based visual inputs. However, the representation capability of this approach is limited by the given categories of the object detection task, leading to information gaps for language understanding [14]. In contrast to the region-based visual embedding, PixelBERT [14] suggests CNN-based visual encoder with random pixel sampling to improve the robustness of visual feature learning and avoid over-fitting [11]. Since there is no applicable off-the-shelf object detector model to extract region-based feature in the medical domain [8], [21], [27], we adopt the CNN-based visual feature embedding.

**3) Downstream Tasks:** VLU and VLG tasks are typical downstream tasks of the VL pre-trained model for tackling more complex tasks that combine vision with language. In this regard, a number of previous works [7], [14], [24], [34] use BERT-based vision-language joint encoder to perform VLU tasks. On the other hand, VLG tasks typically require an encoder for embedding the vision features and a decoder that generates text [33]. Unified VLP [41] conducts these two disparate tasks (VLU and VLG) with a single BERT-based architecture by repeatedly alternating the mask type with a fixed ratio between bidirectional and sequence-to-sequence mask during pre-training. Inspired by this unified pre-training approach, we explore different types of masks and their effects on diverse VLU and VLG downstream tasks.

### III. MATERIALS AND METHODS

#### A. Dataset

We used publicly available MIMIC-CXR [16] and Open-I [9] datasets. MIMIC-CXR [16] contains 377,110 Chest X-ray images and corresponding free-text reports. Also, Open-I dataset contains 3,851 reports and 7,466 Chest X-ray images. Since the dataset contains frontal and lateral view images, it is required to distinguish between view positions [6], [22] to avoid miss-match findings between an image and a report pair. Therefore, given the dominance of the anteroposterior (AP) frontal view in ICU (Intensive Care Units) settings (e.g., 38.89% of all studies containing at least one AP view image), we perform all experiments on unique 91,685 AP view image and associated report pairs following the official split of MIMIC-CXR (train 89,395, valid 759, test 1,531) and 3,547 image-report pairs from the official Open-I dataset. We use Open-I to test the generalization ability of the models, where all models are pre-trained on MIMIC-CXR, then fine-tuned for downstream tasks on a completely unseen Open-I dataset. We pre-process the X-ray image and report data as follows. First, for the X-ray image, we cut out the marginal space of the original image and resize all the images to  $512 \times 512$ , keeping the aspect ratio. Then for the report, we select a

TABLE I  
NOTATION EXPLANATION APPEARING IN THIS SECTION

Notation	Description	Notation	Description
$v$	Chest x-ray image	$\mathbf{p}$	Language position embedding
$\mathbf{v}$	Visual feature	$\mathbf{s}_L$	Language semantic embedding
$\mathbf{l}$	Visual location feature	$\tilde{\mathbf{v}}$	Final visual feature embedding
$\mathbf{s}_V$	Visual semantic embedding	$\tilde{\mathbf{w}}$	Final language feature embedding
$w$	Clinical report	$\tilde{\mathbf{H}}$	Joint embedding
$\mathbf{w}$	Language feature	$\tilde{\mathbf{H}}$	Contextualized embedding

longer description (Findings or Impression section) which may contain detailed information associated with the X-ray imaging.

#### B. VL Pre-Training Model

Our proposed architecture MedViLL is a single BERT-based model that learn unified contextualized vision-language representation. The overall architecture of MedViLL is illustrated in Fig. 2.

**1) Visual Feature Embedding:** We use a CNN to extract visual features from the medical image. The visual features are obtained from the last convolution layer, then flattened along the spatial dimension. Further, we encode the absolute positions of visual input as additional information for explicitly injecting the same body position information in the x-ray images. Given a Chest x-ray image  $v$ , we denote the flattened visual feature obtained from the last CNN layer  $\mathbf{v}$ , and the location feature  $\mathbf{l}$  as follows:

$$\mathbf{v} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_K\}, \quad \mathbf{v}_i \in \mathbb{R}^c \quad (1)$$

$$\mathbf{l} = \{\mathbf{l}_1, \mathbf{l}_2, \dots, \mathbf{l}_K\}, \quad \mathbf{l}_i \in \mathbb{R}^c \quad (2)$$

where  $K$  indicates the number of visual features (i.e., height  $\times$  width) and  $c$  the hidden dimension size (i.e., channel size). The final visual feature embeddings  $\tilde{\mathbf{v}} = \{\tilde{\mathbf{v}}_1, \tilde{\mathbf{v}}_2, \dots, \tilde{\mathbf{v}}_K\}$  are computed as follows:

$$\tilde{\mathbf{v}}_i = \mathbf{v}_i + \mathbf{l}_i + \mathbf{s}_V \quad (3)$$

where  $\mathbf{s}_V$  is a semantic embedding vector shared by all visual feature to differentiate themselves from language embeddings. The final visual features  $\tilde{\mathbf{v}}$  are fed into a fully-connected layer, to be projected into the same embedding space  $\mathbb{R}^d$  as the language embeddings. During pre-training, we randomly sample a subset of the final visual features to avoid overfitting and enhance the semantic knowledge learning of visual input [15]. We use  $k$  to denote the number of sampled visual features, whereas  $K$  denotes the number of all visual features.

**2) Language Feature Embedding:** For language feature embedding, we follow BERT [10] to encode the textual information. A given clinical report  $w$  is first split into a sequence of  $N$  tokens (i.e. subwords)  $\{w_1, \dots, w_N\}$  using the WordPiece tokenizer [38]. The tokens are then converted to vector representations  $\mathbf{w} = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N\}$ ,  $\mathbf{w}_i \in \mathbb{R}^d$  via a lookup table, where  $d$  is the embedding dimension size. We denote position embeddings as  $\mathbf{p} = \{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_N\}$ ,  $\mathbf{p}_i \in \mathbb{R}^d$ . The final language feature embeddings  $\tilde{\mathbf{w}} = \{\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2, \dots, \tilde{\mathbf{w}}_N\}$  are obtained as follows:

$$\tilde{\mathbf{w}}_i = \mathbf{w}_i + \mathbf{p}_i + \mathbf{s}_L \quad (4)$$



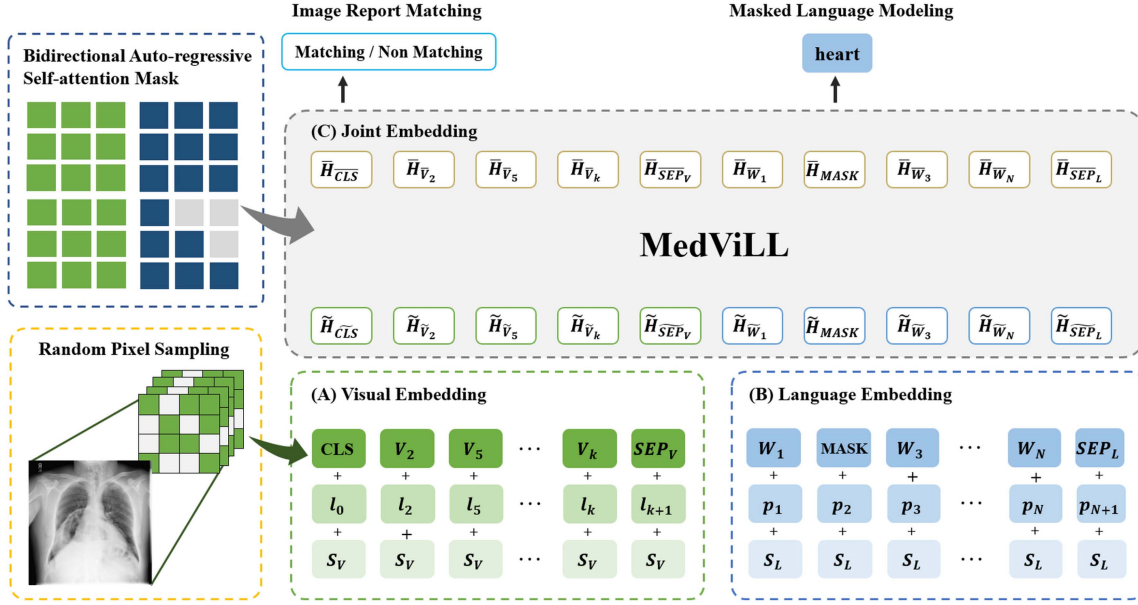


Fig. 2. Architecture of the MedViLL. MedViLL is a single stream BERT model for the cross-modal embedding. Chest X-ray images are randomly sampled from the last feature map of the CNN model as visual inputs. Also, each report is parsed with the BERT tokenizer to get language input. MedViLL is pre-trained with masked language modeling and image report matching tasks, and flexibly applied to VLU and VLG downstream tasks.

where  $s_L$  is a semantic embedding vector shared by all language feature to differentiate themselves from visual embeddings.

**3) Joint Embedding:** After obtaining visual embedding  $\tilde{v} \in \mathbb{R}^d$  and language embeddings  $\tilde{w} \in \mathbb{R}^d$ , we concatenate them to construct the input sequence to the joint embedding component (Fig. 2 (C)). Using additional special tokens CLS and SEP, we define the input to the joint embedding block as  $\tilde{H} = \{\tilde{CLS}, \tilde{v}_1, \dots, \tilde{v}_K, \tilde{SEP}_V, \tilde{w}_1, \dots, \tilde{w}_N, \tilde{SEP}_L\} \in \mathbb{R}^{d \times S}$  where  $S = N + K + 3$ . Note that  $\tilde{CLS}$ ,  $\tilde{SEP}_V$  and  $\tilde{SEP}_L$  are obtained by summing the special tokens with corresponding position and semantic embeddings as in Fig. 2. The contextualized embedding produced by the joint embedding block are denoted as  $\tilde{H} = \{\tilde{CLS}, \tilde{v}_1, \dots, \tilde{v}_K, \tilde{SEP}_V, \tilde{w}_1, \dots, \tilde{w}_N, \tilde{SEP}_L\}$ .

**4) Pre-Training Objectives:** To pre-train MedViLL and align visual features with language features, we take the Masked Language Modeling (MLM) and Image Report Matching (IRM) tasks, which were used in various forms in previous work [7], [14]. For the MLM task, we follow BERT to replace 15% of the input text tokens  $\{w_1, \dots, w_N\}$  with the special MASK token, a random token, or the original token with a probability of 80%, 10% and 10% respectively. The model is trained to recover these masked tokens based on the contextual observation of their surrounding language tokens and the visual tokens, by minimizing the following negative log-likelihood.

$$L_{MLM}(\theta) = -\mathbb{E}_{(v,w) \sim D} [\log P_{\theta}(w_m|v, w_{\setminus m})] \quad (5)$$

where  $\theta$  is the trainable parameters of MedViLL. A pair of images and its corresponding report  $(v, w)$  is sampled from the training set  $D$ , where  $w$  can be divided into the masked tokens  $w_m$  and their complements  $w_{\setminus m}$ .  $\mathbb{E}_{(v,w) \sim D}$  is the average for the training set  $D$ , and  $P_{\theta}(w_m|v, w_{\setminus m})$  is the probability of

$w_m$  given  $v$  and  $w_{\setminus m}$ . IRM task encourages the model to learn both visual and textual features by training the model to predict whether a given pair of image and report  $(v, w)$  is a matching pair or not. During pre-training, we randomly sample both matching image-report pairs and non-matching image-report pairs with 1:1 ratio from the dataset. Note that, however, while selecting a matching pair is straightforward (X-ray images come with a corresponding report), sampling a non-matching pair is not, because two different report can be semantically the same (e.g., “No findings,” “Nothing noticeable.”). Therefore, when sampling for non-matching image-report pairs, we use diagnosis labels. Specifically, in our IRM task, a non-matching report is defined as the ones that are extracted different positive diagnosis labels than the matching report. The joint contextualized embedding  $\tilde{CLS}$  is used to classify whether the input image and report are a matching pair or not, with the following loss function,

$$L_{IRM}(\theta) = -\mathbb{E}_{(v,w) \sim D} [y \log P_{\theta}(v, w)] - \mathbb{E}_{(v,w') \sim D} [(1 - y) \log(1 - P_{\theta}(v, w'))] \quad (6)$$

where  $(v, w)$  denotes a matching image-report pair,  $(v, w')$  a non-matching pair,  $y$  is the label (1 for matching, and 0 for unmatching),  $\mathbb{E}_{(v,w) \sim D}$  is the average for the training set  $D$ , and  $P_{\theta}(v, w)$  is the probability of the  $(v, w)$  being paired.

### C. Self-Attention Mask Schemes

We explore several types of self-attention masks to encourage the model to learn universal multi-modal representations. Bi (Bidirectional) attention mask (Fig. 3(a)) that allows all inputs to interact freely for unconstrained context learning between the visual-language modalities. S2S (Sequence-to-Sequence) causal attention mask (Fig. 3(d)), on the other hand, allows

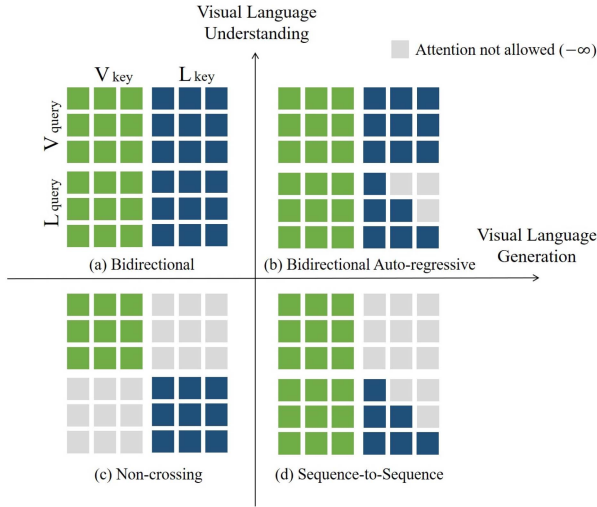


Fig. 3. Self-attention mask schemes. Four types of self-attention masks and the quadrant for the difference in performance in the downstream task of each attention mask. (a) Bidirectional. (b) Bidirectional Auto-regressive. (c) Non-crossing. (d) Sequence-to-sequence self-attention masks.

restricted context learning; language features are only allowed to attend to previous words, while visual features are not allowed to attend to any language features, in order to prevent leaking information from the future. Bi & S2S uses both Bi and S2S masks alternately during pre-training (in every mini-batch, use S2S with 75% chance and Bi with 25% chance) to perform both VLU and VLG downstream tasks. In this work, we propose a new self-attention mask, Bidirectional Auto-Regressive (BAR) (Fig. 3(b)), to closes the gap between Bi and S2S while taking advantage of both. BAR allows image features to be mixed with language features during pre-training (as opposed to S2S mask), while preserving the causal nature of auto-regressive language generation.

The self-attention mask  $M \in \mathbb{R}^{S \times S}$ ,  $S = N + K + 3$  consists of 0s and negative infinities as below.

$$M_{jk} = \begin{cases} 0, & \text{(attention allowed)} \\ -\infty, & \text{(attention not allowed)} \end{cases} \quad j, k = 1, \dots, S. \quad (7)$$

And a single attention head in the self-attention module can be formulated as follows:

$$\text{Attention} = \text{softmax}(SA + M)V, \quad SA = \frac{QK^T}{\sqrt{d_k}} \quad (8)$$

where  $Q, K, V$ , and  $d_k$  indicate queries, keys, values, and dimension of queries and keys respectively [35]. Equation (9) shown at the bottom of this page, where  $q$  and  $k$  indicate query and key vectors respectively. Since the self-attention matrix is computed from the query and key vectors of vision-language modalities according to the (9), the computed self-attention matrix can be divided into 4 subparts of queries and key combinations by modality type.

$$SA_{q,k} = SA_{CLS_q:SEP_{Vq}, CLS_k:SEP_{Vk}} \quad (10)$$

$$+ SA_{CLS_q:SEP_{Vq}, W_{1k}:SEP_{Lk}} \quad (11)$$

$$+ SA_{W_{1q}:SEP_{Lq}, CLS_k:SEP_{Vk}} \quad (12)$$

$$+ SA_{W_{1q}:SEP_{Lq}, W_{1k}:SEP_{Lk}} \quad (13)$$

where (10) is the attention of query and key from vision, (11) is an attention mask of query from the vision and key from language, (12) is an attention mask of query from the language and key from vision, and (13) is an attention mask of query and key from language features. We combine the attention mask matrix  $M$  for the subparts of  $SA$  because adding negative infinity to the calculated attention value will result in zero in the softmax operation.

$$BAR_M = \begin{bmatrix} 0 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 \\ 0 & \dots & -\infty & \dots & -\infty \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 \end{bmatrix} \in \mathbb{R}^{S \times S} \quad (14)$$

Therefore, BAR attention mask (14) allows the attention calculations of all possible combinations except for the (13). Intuitively, this self-attention mask scheme applies auto-regressive attention masks to language modality to enhance joint embedding between vision and language modalities and perform well in both generation and understanding tasks. We implemented four different models each using different types of self-attention masks during pre-training; Bi, S2S, BAR and Bi & S2S. In addition, we also experiment with Non-crossing attention mask (Fig. 3(c)) as a baseline to investigate the impact of multi-modal representation learning. As non-crossing attention mask restricts the interaction between two modalities, we add one additional CLS token at the beginning of the language features, so that both  $\overline{CLS}_V$  and  $\overline{CLS}_L$  can be used for the IRM pre-training task.

$$SA = \begin{bmatrix} CLS_q \cdot CLS_k & \dots & CLS_q \cdot W_{1k} & \dots & CLS_q \cdot SEP_{Lk} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ SEP_{Vq} \cdot CLS_k & \dots & SEP_{Vq} \cdot W_{1k} & \dots & SEP_{Vq} \cdot SEP_{Lk} \\ W_{1q} \cdot CLS_k & \dots & W_{1q} \cdot W_{1k} & \dots & W_{1q} \cdot SEP_{Lk} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ SEP_{Lq} \cdot CLS_k & \dots & SEP_{Lq} \cdot W_{Lk} & \dots & SEP_{Lq} \cdot SEP_{Lk} \end{bmatrix} \quad (9)$$

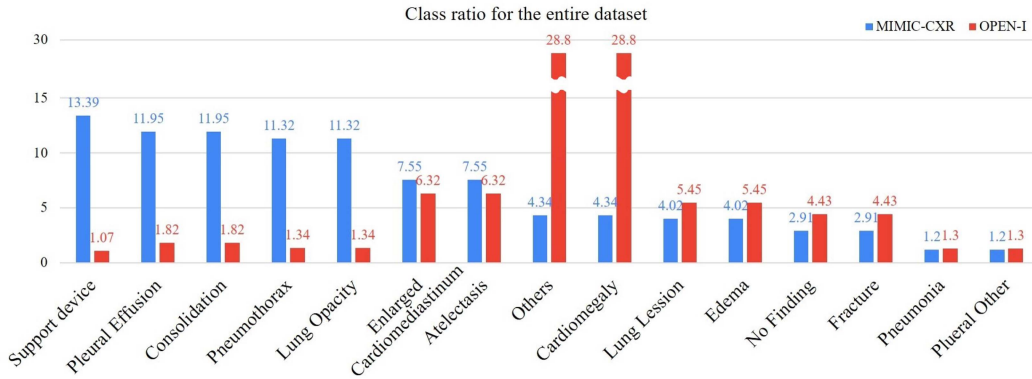


Fig. 4. Dataset Analysis. We compare the distribution of diagnosis labels over the entire dataset. Due to the different scales of the two datasets, each label was represented as a percentage over the entire dataset.

## IV. RESULTS AND DISCUSSION

### A. Dataset Analysis

Although both MIMIC-CXR and Open-I consist of chest X-ray images and report pairs, the two datasets could have different characteristics since they were collected from separate institutions. Specifically, the diagnostic information represented by the two X-ray image sets could be differently distributed. Therefore, to analyze the difference in the distribution of diagnostic labels between two datasets, we compared positive labels acquired from the Chexpert labeler results. As seen in Fig. 4, a mild imbalance was observed in MIMIC-CXR where the class ratios ranged from 13.39% (support devices) to 1.2% (pneumonia, and pleural other). On the other hand, a severe imbalance was observed in Open-I compared to MIMIC-CXR with the maximum class ratios of 28.8% (Others, and cardiomegaly) and the minimum of 1.07% (support devices). This shows that Open-I not only differs from MIMIC-CXR in terms of data volume, but also in terms of clinical properties. Therefore, we believe evaluating the MIMIC-CXR-pre-trained models on Open-I is an appropriate setup to test the generalization capability of the models. The distribution of diagnosis label is illustrated in Fig. 4.

### B. Implementation Details

We use ResNet-50 pre-trained on ImageNet as a visual feature extractor. The input image size is  $(512 \times 512 \times 3)$ , and the last feature map  $(16 \times 16 \times 2048)$  of ResNet-50 is flattened by spatial dimensions and we randomly sample 180 visual features  $(180 \times 2048)$  during pre-training, while we use all features  $(256 \times 2048)$  for every downstream task. To embed text token, each sequence from reports is truncated or padded to 253 tokens in length by considering maximum embedding size. For the joint embedding, we adopt BERT-base architecture which comprised of 12 Transformer layers. Each layer contains 12 attention heads, 768 embedded hidden size and 0.1 drop-out probability. We adopt AdamW optimizer with learning rate  $1e^{-5}$  settings for visual backbone and Transformer. All models were trained on 8 RTX-3090 GPU with the batch size of 128 and 50 epochs for the pre-training model.

### C. Task-Specific Downstream Model Strategy

1) *Diagnosis Classification*: For a given image-report pair, we use the positive labels extracted from the report by the Chexpert labeler as the diagnosis labels. As a single pair could have multiple diagnosis labels up to the maximum of 14 (i.e. multi-label classification), we use 14 linear heads on top of  $\overline{CLS}$  and fine-tune the model using the binary cross-entropy loss. All models are evaluated with the micro average AUROC, and micro average F1 score.

2) *Medical Image-Report Retrieval*: There are two subtasks for medical image-report retrieval, where image-to-report (I2R) retrieval requires the model to retrieve the most relevant report from a large pool of reports given an image, and vice versa for report-to-image (R2I) retrieval. Given an image, any report that contains the same Chexpert diagnosis labels as the original matching report is considered a positive image-report pair, and a negative pair otherwise. The final multi-modal representation  $\overline{CLS}$  is used as the input to a binary classifier to classify the given pair, which is trained by the binary cross-entropy loss. At inference, in each trial, a model is given 100 image-report pairs, and it must use the predicted scores to rank the positive pair as highest as possible. The evaluation metrics are Hit@K, Recall@K, Precision@K ( $K = 5$ ), and mean reciprocal rank (MRR).

3) *Medical Visual Question Answering*: We perform VQA on the VQA-RAD dataset [1], which contains 3,515 question-answer pairs on 315 images (104 head CTs or MRIs, 107 Chest X-rays, and 104 abdominal CTs). As our models are pre-trained on Chest X-ray images, VQA-RAD provides a unique opportunity to study whether the pre-trained models would generalize well beyond the single image domain. Given a pair of an image and a free-text question, we use the final representation  $\overline{CLS}$  to predict a one-hot encoded answer (all possible answers are treated as a single token). The performance was evaluated with accuracy, but separately for the closed questions (i.e. short-form answers such as yes/no) and the open-ended questions (i.e. long-form answers), following the original VQA-RAD paper.

4) *Radiology Report Generation*: The fine-tuning process is same as the MLM pre-training task, except that we fix the

**TABLE II**  
MODEL AUROC AND F1 SCORES FOR THE DIAGNOSIS CLASSIFICATION TASK ON MIMIC-CXR AND OPEN-I

Dataset	Metrics	MedViLL				Baseline					
		Ours	*	R-101	D-121	Bi&S2S	Bi	S2S	Non-crossing	Fine-tuning Only	CNN & Transformer
MIMIC-CXR	avg AUROC	0.980(0.00)	0.926(0.0)	0.918(0.0)	0.962(0.0)	0.979(0.00)	<b>0.984(0.00)</b>	0.982(0.00)	0.980(0.00)	0.969(0.00)	0.831(0.00)
	avg F1	0.839(0.00)	0.699(0.0)	0.666(0.01)	0.771(0.0)	0.846(0.00)	<b>0.852(0.00)</b>	0.846(0.00)	0.824(0.00)	0.807(0.00)	0.491(0.00)
	p-value(avg AUROC)	-	-	-	-	0.005	1.97E-15	0.003	0.254	1.70E-36	3.41E-102
	p-value(avg F1)	-	-	-	-	9.59E-28	7.85E-42	1.62E-26	2.02E-43	4.90E-63	2.70E-122
Open-I	avg AUROC	0.892(0.00)	0.886 (0.00)	0.894 (0.00)	<b>0.897 (0.00)</b>	0.827(0.00)	0.758(0.00)	0.720(0.00)	0.589(0.00)	0.723(0.00)	0.709(0.00)
	avg F1	0.408(0.01)	0.397 (0.00)	0.404 (0.00)	<b>0.409 (0.00)</b>	0.301(0.01)	0.295(0.01)	0.256(0.01)	0.185(0.00)	0.300(0.00)	0.245(0.01)
	p-value(avg AUROC)	-	-	-	-	6.94E-83	4.00E-98	1.23E-101	4.66E-122	1.41E-103	1.35E-109
	p-value(avg F1)	-	-	-	-	1.05E-93	7.04E-95	7.17E-101	2.08E-110	6.49E-94	2.69E-104

\* Indicates the use of all visual features, and R-101 and D-121 indicate the use of ResNet101 and DenseNet121 as visual backbones, respectively. Inference time(ms) on MIMIC-CXR: MedViLL(12.5), Bi&S2S(13), Bi(13), S2S(13), Non-crossing(12.5), Fine-tuning Only(15.5), CNN & Transformer(10.5).

self-attention mask to S2S for all models. At inference, reports can be generated by sequentially recovering the MASK tokens; given visual features followed by a single MASK token, the model can predict the first language token. Then we can replace the first MASK with the sampled token, and a new MASK token is appended. This process is repeated until the model predicts the SEP token as the stop sign. The performance is measured with three metrics: perplexity, clinical efficacy, and BLEU score. Clinical efficacy metrics is obtained by applying the Chexpert labeler on both the original matching report and the generated report. Based on the extracted labels, we can calculate accuracy, precision, recall, and F1 accordingly. Perplexity is used to evaluate the linguistic fluency of the model, while the clinical efficacy is used to evaluate if the model can capture the semantics of the given image. We also report 4-gram BLEU score to evaluate how similar the generated report is to the reference report.

#### D. Downstream Task Result

MedViLL is compared to four pre-training models trained with different attention masks. We also include two more baselines: 1) Fine-tuning Only, which follows the same model architecture as MedViLL, but directly fine-tuned on each downstream task without any pre-training. 2) CNN & Transformer, which uses the CNN module for encoding image only, and the Transformer module (same size as MedViLL) for encoding report only, and the outputs from each module are used for downstream tasks. CNN & Transformer also does not use pre-training. For all tasks, we conduct 30 random multiple-bootstrap experiments and report the mean performance and its standard deviation. Also, we perform a statistical hypothesis test. Based on the average values of the various metrics obtained for each model, we conduct an independent t-test with a significance value of 0.05 to identify the significant pairwise difference of our method against multiple baseline models. In summary, MedViLL achieved the best or second-best performance by analyzing statistical significance with various baseline models in the VLU and VLG tasks. In addition, MedViLL shows superior generalization ability by outperforming most of the models in out-of-domain evaluations.

1) **Diagnosis Classification:** All model performance is shown in Table II. For the Open-I images, we use Chexpert labeler on their MeSH annotations to extract the same set of

diagnosis categories as the MIMIC-CXR dataset. Specifically, Bi and S2S outperform MedViLL with a statistically significant difference in both micro-averaged AUROC and F1 scores, indicating the null hypothesis can be rejected (t-test produced a p-value lower than 0.05). Also, although MedViLL (0.9805) achieves a higher score than Non-crossing (0.9801), it is not statistically meaningful with a p-value of 0.254 in micro-averaged AUROC. However, MedViLL outperforms all other baselines with a statistically meaningful difference in MIMIC-CXR. Moreover, MedViLL outperforms all models statistically significantly when transferred to Open-I. It is also noteworthy that Bi&S2S, which is aimed to take advantage of both the bidirectional mask and the S2S mask demonstrates much better generalization capability compared to the two individual masks.

2) **Medical Image-Report Retrieval:** Table III shows the performance of Image-to-Report and Report-to-Image retrieval. We report the p-value of MRR for the performance of both tasks since MRR is a rank-aware evaluation metrics compared to other metrics. We can observe that all pre-trained models significantly outperform naive baselines (Fine-tune Only and CNN & Transformer) for the MIMIC-CXR dataset. In Report-to-Image retrieval, while MedViLL achieves lower performance than Bi and S2S with a statistically significant difference in MIMIC-CXR, MedViLL outperforms all baselines when fine-tuned on the unseen Open-I dataset except for Bi. However, although Bi outperforms MedViLL in Open-I, there is no statistically significant difference between both models with a p-value of 0.843. In Image-to-Report retrieval, S2S statistically outperforms MedViLL in MIMIC-CXR, but MedViLL is superior to all baselines with a statistically meaningful difference in Open-I. It is notable that the naive baselines, while severely underperforming for MIMIC-CXR, show substantially increased performance for Open-I. We believe this is due to the Open-I being a significantly smaller dataset than MIMIC-CXR, with only two Chexpert labels (Others, and Cardiomegaly) mostly dominating the label space.

3) **Medical Visual Question Answering:** Table IV shows the VQA accuracy when models were fine-tuned with all image types ('ALL'), and with only the Chest X-ray images ('CHEST'). We can see that MedViLL significantly outperforms MEVF [25], the state-of-the-art model for the VQA-RAD dataset, indicating the effectiveness of the multi-modal pre-training for this complex multi-modal reasoning task. We can also see that MedViLL shows comparable performance as the bidirectional



TABLE III  
MEDICAL IMAGE-REPORT RETRIEVAL PERFORMANCE ON MIMIC-CXR AND OPEN-I

Task	Models	MIMIC-CXR					OpenI				
		MRR	H@5	R@5	P@5	p-value	MRR	H@5	R@5	P@5	p-value
Report-to-Image	MedViLL(Ours)	56.5(0.01)	77.0(0.01)	47.4(0.01)	19.9(0.00)	-	51.3(0.01)	73.0(0.01)	12.9(0.00)	31.7(0.00)	-
	MedViLL*	49.0(0.1)	69.7(0.1)	41.9(0.1)	17.6(0.1)	-	54.8(0.1)	74.7(0.1)	13.4(0.1)	35.6(0.1)	-
	MedViLL(R-101)	47.5(0.1)	69.3(0.1)	41.5(0.1)	17.4(0.1)	-	53.1(0.1)	73.2(0.1)	12.8(0.1)	35.8(0.1)	-
	MedViLL(D-121)	49.3(0.1)	70.0(0.1)	42.2(0.1)	17.75(0.1)	-	<b>56.4(0.1)</b>	74.5(0.1)	<b>14.0(0.1)</b>	<b>36.7(0.1)</b>	-
	Bi&S2S	55.5(0.01)	76.7(0.01)	46.7(0.01)	19.7(0.00)	1.20E-05	46.4(0.01)	68.1(0.01)	10.5(0.00)	28.8(0.01)	3.71E-27
	Bi	58.0(0.01)	78.2(0.01)	48.2(0.01)	20.2(0.00)	1.60E-10	51.4(0.01)	<b>74.8(0.01)</b>	13.3(0.00)	32.0(0.01)	0.843
	S2S	<b>58.8(0.01)</b>	<b>79.1(0.01)</b>	<b>48.9(0.01)</b>	<b>20.3(0.00)</b>	1.89E-18	48.6(0.01)	67.2(0.01)	10.3(0.01)	32.9(0.01)	2.28E-14
	Non-crossing	54.7(0.01)	77.0(0.01)	47.2(0.01)	19.5(0.00)	4.07E-12	48.6(0.01)	68.4(0.01)	11.2(0.00)	31.1(0.01)	3.88E-18
	Fine-tuning Only	41.8(0.01)	61.6(0.01)	35.8(0.01)	15.8(0.00)	3.14E-53	36.9(0.01)	54.4(0.01)	5.4(0.00)	20.7(0.01)	1.72E-53
Image-to-Report	CNN & Transformer	11.4(0.01)	15.2(0.02)	5.1(0.00)	3.6(0.01)	9.43E-71	36.2(0.04)	56.6(0.04)	5.0(0.00)	21.4(0.04)	4.94E-19
	MedViLL(Ours)	55.8 (0.01)	75.5(0.01)	47.1(0.01)	19.7(0.00)	-	<b>50.4(0.01)</b>	63.8(0.01)	12.9(0.00)	35.5(0.01)	-
	MedViLL*	48.6(0.1)	68.0(0.1)	42.4(0.1)	17.4(0.1)	-	48.5(0.1)	62.2(0.1)	12.5(0.1)	34.3(0.1)	-
	MedViLL(R-101)	46.3(0.1)	65.5(0.1)	39.8(0.1)	17.1(0.1)	-	47.2(0.1)	62.0(0.1)	11.8(0.1)	32.5(0.1)	-
	MedViLL(D-121)	48.6(0.1)	69.4(0.1)	42.7(0.1)	18.38(0.1)	-	46.5(0.1)	60.9(0.1)	12.8(0.1)	31.7(0.1)	-
	Bi&S2S	54.5(0.01)	75.5(0.01)	47.8(0.01)	19.9(0.00)	6.32E-08	45.8(0.01)	54.0(0.01)	10.1(0.00)	35.8(0.00)	8.55E-29
	Bi	56.7(0.01)	76.3(0.01)	47.6(0.01)	20.2(0.00)	0.0002	48.5(0.01)	<b>65.8(0.01)</b>	<b>13.7(0.00)</b>	32.3(0.01)	3.17E-12
	S2S	<b>57.9(0.01)</b>	<b>78.5(0.01)</b>	<b>49.7(0.01)</b>	<b>20.7(0.00)</b>	2.72E-13	45.4(0.01)	53.6(0.01)	8.9(0.00)	<b>36.9(0.00)</b>	6.84E-31
	Non-crossing	54.6(0.01)	75.7(0.01)	47.6(0.01)	20.0(0.00)	3.84E-07	42.6(0.01)	61.2(0.01)	11.0(0.00)	28.0(0.01)	1.15E-40
	Fine-tuning Only	41.4(0.01)	60.8(0.01)	36.3(0.01)	15.7(0.00)	5.56E-56	45.2(0.00)	49.7(0.01)	5.1(0.00)	35.0(0.00)	2.64E-29
	CNN & Transformer	12.0(0.02)	15.3(0.02)	5.1(0.00)	4.0(0.01)	1.09E-52	37.9(0.06)	54.0(0.06)	5.0(0.00)	23.0(0.06)	1.16E-12

MedViLL\* indicates the use of all visual features, and R-101 and D-121 indicate the use of ResNet101 and DenseNet121 as visual backbones, respectively. Inference time(ms) of Report-to-Image and Image-to-Report on MIMIC-CXR: MedViLL(7.6, 7.8), Bi&S2S(8.2, 7.8), Bi(7.6, 7.6), S2S(7.6, 7.7), Non-crossing(7.7, 7.8), Fine-tuning Only(7.8, 7.6), CNN & Transformer(5.3, 5.3).

TABLE IV  
MODEL ACCURACY ON THE VQA-RAD DATASET

Models	ALL				CHEST			
	O.E.	C.E.	p-value of O.E.	p-value of C.E.	O.E.	C.E.	p-value of O.E.	p-value of C.E.
MedViLL(Ours)	<b>0.595(0.032)</b>	0.777(0.071)	-	-	0.587(0.033)	<b>0.782(0.123)</b>	-	-
MedViLL*	0.512(0.012)	0.743(0.009)	-	-	0.572(0.021)	0.736(0.006)	-	-
MedViLL(R-101)	0.548(0.019)	0.781(0.002)	-	-	0.602(0.012)	0.773(0.022)	-	-
MedViLL(D-121)	0.593(0.001)	<b>0.787(0.004)</b>	-	-	<b>0.612(0.01)</b>	0.781(0.003)	-	-
Bi&S2S	0.541(0.038)	0.76(0.027)	2.93E-07	0.224	0.566(0.074)	0.766(0.035)	0.164	0.519
Bi	0.58(0.038)	0.784(0.03)	0.124	0.643	0.562(0.04)	0.767(0.035)	0.013	0.549
S2S	0.505(0.042)	0.73(0.025)	1.81E-12	0.002	0.517(0.07)	0.723(0.048)	1.57E-05	0.021
Non-crossing	0.531(0.015)	0.734(0.017)	5.58E-12	0.003	0.474(0.083)	0.732(0.03)	3.94E-08	0.043
Fine-tuning Only	0.232(0.019)	0.649(0.026)	2.98E-43	5.38E-11	0.124(0.014)	0.606(0.035)	1.08E-42	1.50E-08
CNN & Transformer	0.24(0.029)	0.667(0.015)	7.66E-46	2.70E-9	0.124(0.067)	0.523(0.033)	7.60E-32	1.62E-12
MEVF [25]	0.407	0.741	-	-	-	-	-	-

O.E. stands for Open-ended question and C.E. stands for close-ended question. MedViLL\* indicates the use of all visual features, and R-101 and D-121 indicate the use of ResNet101 and DenseNet121 as visual backbones, respectively. For MEVF [?], we used the reported results from the original paper. Inference time(ms) on MIMIC-CXR: MedViLL(19.46), Bi&S2S(19.52), Bi(19.43), S2S(19.51), Non-crossing(19.58), Fine-tuning Only(19.61), CNN & Transformer(17.42).

mask which allows unrestricted interaction between all text and vision features. In a statistical analysis, MedViLL shows significantly higher performance than all models for O.E. (open-ended questions) of ‘ALL’ and ‘CHEST’. However, there was no statistically significant difference for Bi&S2S of ‘CHEST’ and Bi of ‘ALL,’ (p-values of 0.164, and 0.124, respectively). For C.E. (close-ended questions) of ‘ALL,’ Bi performed the best of all the models, followed by MedViLL. However, this result is not statistically meaningful, obtaining a p-value of 0.643 between Bi and MedViLL. Also, for C.E. of ‘CHEST,’ MedViLL outperforms all baselines, but it is not statistically significant against Bi&S2S (p-value of 0.519) and Bi (p-value of 0.549).

4) **Radiology Report Generation:** Table V shows the report generation performance of all models. For this task, we also implemented TieNet [37] as a baseline, which is a widely used CNN-RNN based attention model for report generation. We can see that the models pre-trained with auto-regressive

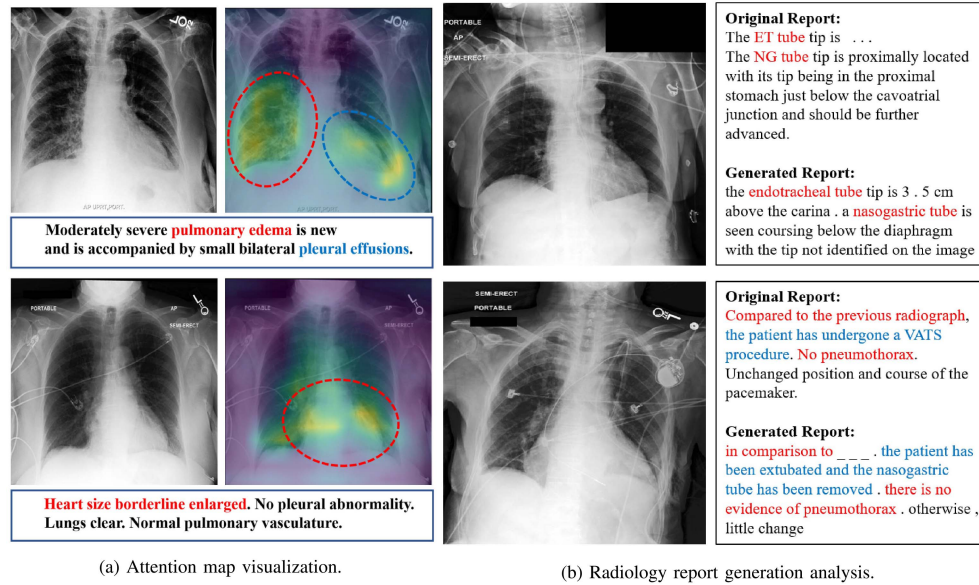
manners (MedViLL, Bi&S2S, S2S) all significantly outperform the other models in terms of both perplexity (except for TieNet) and clinical efficacy metrics for the MIMIC-CXR dataset. Among the clinical efficacy metrics, we report the p-value of F1 score because the F1 score is a balanced measure of both precision and recall and allows us to better capture the true performance here in light of the strong class imbalance of dataset. MedViLL achieved the best performance on the MIMIC dataset with a statistically significant difference. MedViLL seems to best capture the semantics embedded in the image given the highest F1 score. When fine-tuned on the unseen Open-I dataset, TieNet performed the best of all the models, followed by MedViLL. However, this result is not statistically significant between TieNet and MedViLL with a p-value of 0.2181, indicating that the performance of the two models is similar. As opposed to its generally favorable performance in VLU tasks, the bidirectional mask seems to



**TABLE V**  
REPORT GENERATION PERFORMANCE IN TERMS OF PERPLEXITY AND LABEL ACCURACY, PRECISION RECALL AND F1 AND BLEU4

Dataset	Models	Perplexity ( $\downarrow$ )	Accuracy ( $\uparrow$ )	Precision ( $\uparrow$ )	Recall ( $\uparrow$ )	F1 Score ( $\uparrow$ )	BLEU4 ( $\uparrow$ )	p-value
MIMIC-CXR	MedViLL(Ours)	4.185(0.022)	<b>0.841(0.003)</b>	<b>0.698(0.002)</b>	<b>0.559(0.004)</b>	<b>0.621(0.002)</b>	0.066(0.001)	-
	MedViLL*	5.04(0.01)	0.780(0.016)	0.643(0.038)	0.417(0.026)	0.505(0.027)	0.072(0.009)	-
	MedViLL(R-101)	<b>3.91(0.003)</b>	0.833(0.02)	0.652(0.004)	0.472(0.024)	0.526(0.003)	0.116(0.02)	-
	MedViLL(D-121)	4.02(0.002)	0.81(0.003)	0.628(0.02)	0.425(0.004)	0.484(0.005)	<b>0.126(0.003)</b>	-
	Bi&S2S	6.515(0.12)	0.786(0.007)	0.619(0.003)	0.435(0.009)	0.511(0.006)	0.066(0.001)	4.17E-43
	Bi	849.67(5.225)	0.637(0.004)	0.283(0.007)	0.07(0.024)	0.11(0.032)	0.015(0.004)	1.11E-36
	S2S	4.258(0.069)	0.797(0.007)	0.662(0.004)	0.448(0.01)	0.534(0.007)	0.043(0.001)	2.32E-38
	Non-crossing	718.122(9.484)	0.634(0.005)	0.277(0.013)	0.076(0.004)	0.12(0.005)	0.007(0.001)	2.30E-75
Open-I	Fine-tuning Only	224.343(0.204)	0.664(0.003)	0.417(0.012)	0.305(0.006)	0.352(0.005)	0.009(0.004)	4.14E-65
	TieNet	4.132(0.033)	0.687(0.003)	0.487(0.003)	0.380(0.006)	0.426(0.006)	0.123(0.002)	7.17E-54
	MedViLL(Ours)	5.637(0.259)	0.734(0.001)	0.512(0.002)	0.594(0.001)	0.55(0.001)	0.049(0.001)	-
	MedViLL*	5.06(0.081)	0.790(0.018)	0.577(0.015)	0.294(0.089)	0.382(0.081)	0.044(0.010)	-
	MedViLL(R-101)	<b>3.815(0.002)</b>	0.812(0.003)	<b>0.675(0.006)</b>	0.367(0.003)	0.466(0.001)	0.075(0.031)	-
	MedViLL(D-121)	3.998(0.013)	<b>0.818(0.011)</b>	0.582(0.03)	0.385(0.012)	0.472(0.002)	0.071(0.009)	-
	Bi&S2S	15.97(1.071)	0.712(0.003)	0.497(0.003)	0.369(0.006)	0.423(0.004)	0.024(0.01)	4.02E-52
	Bi	787.66(55.492)	0.686(0.004)	0.356(0.025)	0.103(0.006)	0.16(0.008)	0.015(0.004)	2.14E-52
	S2S	4.732(0.537)	0.736(0.003)	0.517(0.002)	0.538(0.004)	0.527(0.002)	0.043(0.002)	1.76E-44
	Non-crossing	217.27(12.139)	0.693(0.003)	0.337(0.025)	0.085(0.005)	0.135(0.007)	0.002(0.001)	1.46E-57
	Fine-tuning Only	292.60(19.858)	0.684(0.003)	0.291(0.023)	0.073(0.035)	0.112(0.047)	0.006(0.002)	7.02E-30
	TieNet	7.901(0.483)	0.732(0.007)	0.517(0.013)	<b>0.610(0.017)</b>	<b>0.553(0.013)</b>	<b>0.189(0.005)</b>	0.2181

MedViLL\* indicates the use of all visual features, and R-101 and D-121 indicate the use of ResNet101 and DenseNet121 as visual backbones, respectively. Inference time(ms) on MIMIC-CXR: MedViLL(32.81), Bi&S2S(33.55), Bi(32.54), S2S(33.04), Non-crossing(32.71), Fine-tuning Only(32.92).



**Fig. 5.** Qualitative results and analysis. We visualize the attention regions extracted from the MedViLL (a). Also, we compare the generated report with the original report on the same chest X-ray image (b).

be evidently harmful for the VLG task, most likely due to its incompatibility with the auto-regressive nature of VLG. Interestingly, we found N-gram based measures to be a suboptimal measure for report generation by showing that all models except TieNet achieved low BLEU scores; where the original report contains abbreviated terms (e.g. “ET tube”) models would generate expanded terms (e.g. “endotracheal tube”) and vice versa.

**5) Ablation Studies of Visual Features:** To demonstrate the effectiveness of visual features, we performed ablation experiments with sampling methods (sampling all visual features) and

various backbone (ResNet-101 [11] and DenseNet-121 [13]). The detailed training method for each model is the same as that of MedViLL. As shown in Table II, Table III, Table IV, and Table V, all results reached poor performance in both in-domain and out-of-domain evaluation when using all the visual features of the CNN. These results show the difficulty of reaching good performance by over-fitting the MIMIC-CXR training set. Also, there are performance differences for each downstream task when using various backbones. We believe that better performance can be obtained with various experiments on the model architecture or hyper-parameter tuning.

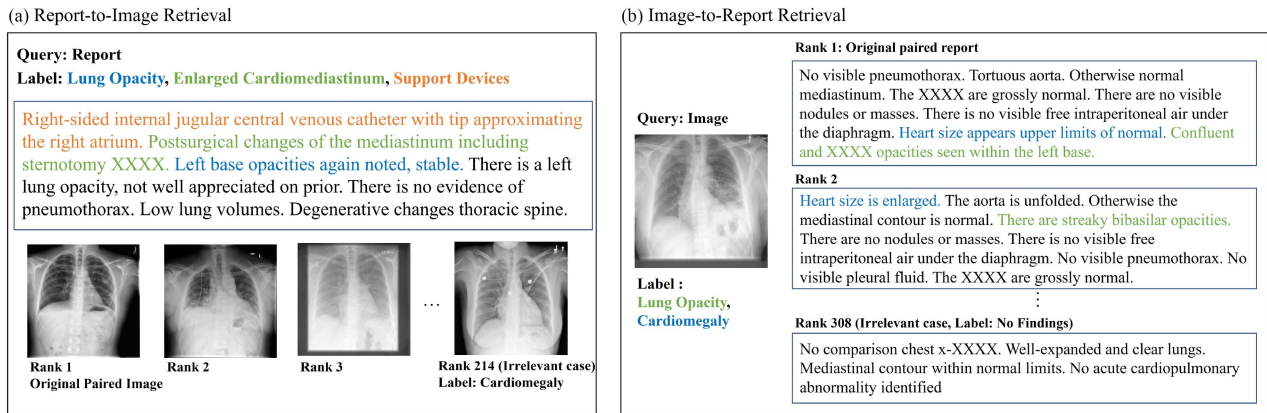


Fig. 6. Case study of Medical Image-Report Retrieval. Given a report or image as a query, we show the retrieved images or reports in (a) and (b) respectively. In (a), the given report is annotated by the Chexpert labeler with the diagnosis labels Lung Opacity (blue), Enlarged Cardiomediastinum (green), and Support Devices (orange). The top three images also contain the same labels as the given query while the last image (Rank 214) is irrelevant to the given query. In (b), the top three reports also contain the same labels recognized by the Chexpert labeler, Lung Opacity (green) and Cardiomegaly (blue). Note that the last report (Rank 308) does not contain any label, and is irrelevant to the given query image.

## E. Qualitative Results and Analysis

1) *Attention Map Visualization*: We visualize the attention maps of the intermediate Transformer layers of MedViLL in order to gain qualitative insight of the cross-modality alignment between text tokens and image features as shown in Fig. 5(a). Although MedViLL uses only report and images for training without any annotations, it can well attend to the disease-discriminatory regions written in the reports, as confirmed by a professional cardiologist. This suggests the potential of MedViLL's capability to explain its learned representations to human users for smoother real-world deployment.

2) *Radiology Report Generation*: We compare the generated reports with the original report. As confirmed by a professional cardiologist, Fig. 5(b) shows that MedViLL is able to generate clinically appropriate reports. Specifically, the blue text in the original report describes the completion of VATS (i.e. video-assisted thoracic surgery). Interestingly, the generated report describes extubation and nasogastric tube removal, which is a part of VATS. This indicates that the BLEU score is not an appropriate measure to evaluate report generation especially in the medical domain.

3) *Medical Image-Report Retrieval*: We retrieved the cases pooled from 1,536 studies in the test set (Fig. 6). As confirmed by a cardiologist, Fig. 6 demonstrate the clinical understanding of MedViLL. Specifically, we can observe that the results in the top-3 retrieved samples all share the same diagnosis labels as the given query; all top three images in Fig. 6(a) are labeled with "Lung Opacity", "Enlarged Cardiomediastinum" and "Support Devices" as the query report, and all top three reports in Fig. 6(b) contain the same labels "Cardiomegaly" and "Lung Opacity" as the query image. Note that the samples in the low rank contain labels irrelevant to the given query.

## V. CONCLUSION

In this study, we propose a multi-modal pre-training model MedViLL, which uses a novel self-attention scheme to

flexibly adapt to multiple downstream tasks of vision-language understanding and generation. By statistically and rigorously evaluating MedViLL on all four downstream tasks with three radiographic image-report datasets, we empirically demonstrated the superior performance of MedViLL against various baselines including task-specific architectures. Despite the impressive performance of MedViLL, this is just the beginning of the vision-language representation learning in the medical domain, and we plan to expand this approach to more diverse settings such as multi-view Chest X-ray studies or a sequence of studies over time.

## REFERENCES

- [1] R. Alizadehsani et al., "A database for using machine learning and data mining techniques for coronary artery disease diagnosis," *Sci. Data*, vol. 6, no. 1, pp. 1–13, 2019.
- [2] E. Alsentzer et al., "Publicly available clinical bert embeddings," in *Proc. 2nd Clin. Natural Lang. Process. Workshop*, 2019, pp. 72–78.
- [3] P. Anderson et al., "Bottom-up and top-down attention for image captioning and visual question answering," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 6077–6086.
- [4] S. Biswal, P. Zhuang, A. Pyrros, N. Siddiqui, S. Koyejo, and J. Sun, "Emixer: End-to-end multimodal x-ray generation via self-supervision," in *Proc. Mach. Learn. Healthcare Conf.*, to be published.
- [5] E. Bugliarello, R. Cotterell, N. Okazaki, and D. Elliott, "Multimodal pretraining unmasked: A meta-analysis and a unified framework of vision-and-language BERTs," *Trans. Assoc. Comput. Linguistics*, vol. 9, pp. 978–994, Sep. 2021.
- [6] W. W. Chapman, M. Fizman, B. E. Chapman, and P. J. Haug, "A comparison of classification algorithms to automatically identify chest x-ray reports that support pneumonia," *J. Biomed. Inform.*, vol. 34, no. 1, pp. 4–14, 2001.
- [7] Y.-C. Chen et al., "Uniter: Universal image-text representations learning," in *Proc. Eur. Conf. Comput. Vis.*, 2020, pp. 104–120.
- [8] V. Cheplygina, M. de Bruijne, and J. P. Pluim, "Not-so-supervised: A survey of semi-supervised, multi-instance, and transfer learning in medical image analysis," *Med. Image Anal.*, vol. 54, pp. 280–296, 2019.
- [9] D. Demner-Fushman et al., "Preparing a collection of radiology examinations for distribution and retrieval," *J. Amer. Med. Inform. Assoc.*, vol. 23, no. 2, pp. 304–310, 2016.
- [10] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. Conf. North Amer. Chapter Assoc. Computat. Linguistics: Hum. Lang. Technol.*, 2019, pp. 4171–4186.

- [11] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [12] T.-M. H. Hsu, W.-H. Weng, W. Boag, M. McDermott, and P. Szolovits, "Unsupervised multimodal representation learning across medical images and reports," 2018, *arXiv:1811.08615*.
- [13] G. Huang, Z. Liu, L. V. D. Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 4700–4708.
- [14] Z. Huang, Z. Zeng, B. Liu, D. Fu, and J. Fu, "Pixel-Bert: Aligning image pixels with text by deep multi-modal transformers," 2020, *arXiv:2004.00849*.
- [15] J. Irvin et al., "Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison," in *Proc. AAAI Conf. Artif. Intell.*, 2019, vol. 33, pp. 590–597.
- [16] A. E. Johnson et al., "MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports," *Sci. Data*, vol. 6, no. 1, pp. 1–8, 2019.
- [17] C. E. Kahn Jr. et al., "Toward best practices in radiology reporting," *Radiology*, vol. 252, no. 3, pp. 852–856, 2009.
- [18] J. Kalpathy-Cramer and W. Hersh, "Multimodal medical image retrieval: Image categorization to improve search precision," in *Proc. Int. Conf. Multimedia Inf. Retrieval*, 2010, pp. 165–174.
- [19] J. Lee et al., "BioBERT: A pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020.
- [20] Y. Li, H. Wang, and Y. Luo, "A comparison of pre-trained vision-and-language models for multimodal representation learning across medical images and reports," in *Proc. IEEE Int. Conf. Bioinf. Biomed.*, 2020, pp. 1999–2004.
- [21] G. Litjens et al., "A survey on deep learning in medical image analysis," *Med. Image Anal.*, vol. 42, pp. 60–88, 2017.
- [22] G. Liu et al., "Clinically accurate chest X-ray report generation," in *Proc. Mach. Learn. Healthcare Conf.*, 2019, pp. 249–269.
- [23] H. Liu et al., "Use of bert (bidirectional encoder representations from transformers)-based deep learning method for extracting evidences in chinese radiology reports: Development of a computer-aided liver cancer diagnosis framework," *J. Med. Internet Res.*, vol. 23, no. 1, 2021, Art. no. e19689.
- [24] J. Lu, D. Batra, D. Parikh, and S. Lee, "ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2019, pp. 13–23.
- [25] B. D. Nguyen, T.-T. Do, B. X. Nguyen, T. Do, E. Tjiputra, and Q. D. Tran, "Overcoming data limitation in medical visual question answering," in *Proc. Int. Conf. Med. Image Comput. Comput.-Assist. Interv.*, 2019, pp. 522–530.
- [26] D. Qi, L. Su, J. Song, E. Cui, T. Bharti, and A. Sacheti, "ImageBERT: Cross-modal pre-training with large-scale weak-supervised image-text data," 2020, *arXiv:2001.07966*.
- [27] M. Raghu, C. Zhang, J. Kleinberg, and S. Bengio, "Transfusion: Understanding transfer learning for medical imaging," in *Proc. Adv. Neural Inf. Process. Syst.*, 2019, pp. 3347–3357.
- [28] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 6, pp. 1137–1149, Jun. 2017.
- [29] T. Schlegl, S. M. Waldstein, W.-D. Vogl, U. Schmidt-Erfurth, and G. Langs, "Predicting semantic descriptions from medical images with convolutional neural networks," in *Proc. Int. Conf. Inf. Process. Med. Imag.*, 2015, pp. 437–448.
- [30] L. H. Schwartz, D. M. Panicek, A. R. Berk, Y. Li, and H. Hricak, "Improving communication of diagnostic radiology findings through structured reporting," *Radiology*, vol. 260, no. 1, pp. 174–181, 2011.
- [31] A. Smit, S. Jain, P. Rajpurkar, A. Pareek, A. Ng, and M. Lungren, "Combining automatic labelers and expert annotations for accurate radiology report labeling using BERT," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, Nov. 2020, pp. 1500–1519.
- [32] C. Sun, F. Baradel, K. Murphy, and C. Schmid, "Learning video representations using contrastive bidirectional transformer," 2019, *arXiv:1906.05743*.
- [33] C. Sun, A. Myers, C. Vondrick, K. Murphy, and C. Schmid, "VideoBERT: A joint model for video and language representation learning," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2019, pp. 7464–7473.
- [34] H. Tan and M. Bansal, "Lxmert: Learning cross-modality encoder representations from transformers," in *Proc. Conf. Empirical Methods Natural Lang. Process. 9th Int. Joint Conf. Natural Lang. Process.*, 2019, pp. 5100–5111.
- [35] A. Vaswani et al., "Attention is all you need," in *Proc. Adv. Neural Inf. Process. Syst.*, 2017, pp. 5998–6008.
- [36] 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 *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 2097–2106.
- [37] X. Wang, Y. Peng, L. Lu, Z. Lu, and R. M. Summers, "TieNet: Text-image embedding network for common thorax disease classification and reporting in chest X-rays," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 9049–9058.
- [38] Y. Wu et al., "Google's neural machine translation system: Bridging the gap between human and machine translation," 2016, *arXiv:1609.08144*.
- [39] Q. Xia et al., "XGPT: Cross-modal generative pre-training for image captioning," in *Proc. CCF Int. Conf. Natural Lang. Process. Chin. Comput.*, 2021, pp. 786–797.
- [40] X. Yang, M. Ye, Q. You, and F. Ma, "Writing by memorizing: Hierarchical retrieval-based medical report generation," in *Proc. 59th Annu. Meeting Assoc. Computat. Linguistics 11th Int. Joint Conf. Natural Lang. Process.*, 2021, pp. 5000–5009.
- [41] L. Zhou, H. Palangi, L. Zhang, H. Hu, J. J. Corso, and J. Gao, "Unified vision-language pre-training for image captioning and VQA," in *Proc. AAAI Conf. Artif. Intell.*, 2020, pp. 13041–13049.