# Towards Unifying Medical Vision-and-Language Pre-training via Soft Prompts

**Zhihong Chen**[1,2*]     **Shizhe Diao**[3*]

**Benyou Wang**[1,2†]     **Guanbin Li**[4†]     **Xiang Wan**[2]

[1]The Chinese University of Hong Kong, Shenzhen   [2]Shenzhen Research Institute of Big Data
[3]The Hong Kong University of Science and Technology   [4]Sun Yat-sen University
zhihongchen@link.cuhk.edu.cn   sdiaoaa@connect.ust.hk
wangbenyou@cuhk.edu.cn   liguanbin@mail.sysu.edu.cn   wanxiang@sribd.com

## Abstract

Medical vision-and-language pre-training (Med-VLP) has shown promising improvements on many downstream medical tasks owing to its applicability to extracting generic representations from medical images and texts. Practically, there exist two typical types, *i.e.*, the fusion-encoder type and the dual-encoder type, depending on whether a heavy fusion module is used. The former is superior at multi-modal tasks owing to the sufficient interaction between modalities; the latter is good at uni-modal and cross-modal tasks due to the single-modality encoding ability. To take advantage of these two types, we propose an effective yet straightforward scheme named PTUnifier to unify the two types. We first unify the input format by introducing visual and textual prompts, which serve as a feature bank that stores the most representative images/texts. By doing so, a single model could serve as a *foundation model* that processes various tasks adopting different input formats (*i.e.*, image-only, text-only, and image-text-pair). Furthermore, we construct a prompt pool (instead of static ones) to improve diversity and scalability. Experimental results show that our approach achieves state-of-the-art results on a broad range of tasks, spanning uni-modal tasks (*i.e.*, image/text classification and text summarization), cross-modal tasks (*i.e.*, image-to-text generation and image-text/text-image retrieval), and multi-modal tasks (*i.e.*, visual question answering), demonstrating the effectiveness of our approach. Note that the adoption of prompts is orthogonal to most existing Med-VLP approaches and could be a beneficial and complementary extension to these approaches.[1]

## 1 Introduction

Medical data is multi-modal in general, among which vision and language are two critical modalities. It includes visual data (*e.g.*, radiography, magnetic resonance imaging, and computed tomography) and textual data (*e.g.*, radiology reports, and medical texts). More importantly, such images and texts are pair-collected in routine clinical practice (*e.g.*, X-ray images and their corresponding radiology reports). Medical vision-and-language pre-training (Med-VLP) aims to learn generic representation from large-scale medical image-text pairs and then transfer it to various medical tasks, which is believed to be beneficial in addressing the data scarcity problem in the medical field.

---

[*]Equal contributions.
[†]Corresponding authors.
[1]Work in progress. The code will be released at https://github.com/zhjohnchan/PTUnifier.

Recently, substantial progress has been made toward research on Med-VLP [66, 31, 19, 42, 7]. In general, most existing Med-VLP models can be classified into two types: the dual-encoder type and the fusion-encoder type, where the former encodes images and texts separately to learn uni-modal/cross-modal representations following a shallow interaction layer (*i.e.*, an image-text contrastive layer), and the latter performs an early fusion of the two modalities through the self-attention/co-attention mechanisms to learn multi-modal representations.[2] For dual-encoders, the purpose of existing studies [66, 19, 44, 60, 57, 61, 3] is to develop label-efficient algorithms to learn effective uni-modal/cross-modal representations since large-scale manually labeled datasets are difficult and expensive to obtain for medical images. The learned representations can improve the *effectiveness* of uni-modal (*i.e.*, vision-only or language-only) tasks[3] and the *efficiency* of cross-modal (*i.e.*, image-to-text or text-to-image) retrieval tasks significantly. For fusion-encoders, existing studies [31, 24, 42, 7, 8] aim to jointly process these two modalities with an early interaction to learn multi-modal representations to solve those tasks requiring multi-modal reasoning (*e.g.*, medical visual question answering and medical image-text classification). However, it seems that "*you can't have your cake and eat it, too.*": the fusion-encoders can not perform uni-modal tasks effectively and cross-modal tasks efficiently due to the lack of single-modal encoding, while the dual-encoders underperform on multi-modal tasks owing to the insufficient interaction between modalities as shown in Figure 1(a).

In this paper, we aim to learn a unified medical vision-and-language pre-trained model. Although there exist some solutions [4, 52] to achieve a similar goal in the general domain, we propose an architecture- and task-agnostic approach named PTUnifier, which is much simpler and lighter-weighted. Technically, we develop the designs from the following perspectives: (i) *Compatibility*: we introduce visual and textual prmopts to make the Med-VLP model compatible with different kinds of inputs (*i.e.*, image-only inputs, text-only inputs, and image-text pairs); (ii) *Scalability*: we improve the diversity of the prompts by constructing prompt pools for different modalities from which different inputs are able to select their corresponding prompts, which enhances the capacity and makes it scalable to larger-scale Med-VLP. As a result, the proposed approach can be employed in unifying Med-VLP with many existing VLP model architectures (*e.g.*, classic one [28] or even a single vanilla Transformer model) and does not require extra modality-dependent architectures, resulting in better applicability. We perform the pre-training on three large-scale medical image-text datasets, *i.e.*, ROCO [47], MedICaT [54], and MIMIC-CXR [23]. To verify the effectiveness of our approach and facilitate further research, we construct a medical vision-language benchmark including uni-modal tasks (*i.e.*, image classification (IC) for vision and text classification (TC) and text summarization (TS) for language), cross-modal tasks (*i.e.*, image-to-text retrieval (ITR), text-to-image retrieval (TIR), and image-to-text generation[4] (ITG)), and multi-modal tasks (*i.e.*, visual question answering (VQA)). The proposed PTUnifier achieves excellent performance on all datasets, demonstrating its effectiveness.

## 2   Related Work

**Vision-and-Language Pre-training (VLP)**   Motivated by the success of the self-supervised pre-training recipe in natural language processing (NLP) (*e.g.*, BERT [13]) and computer vision (CV) (*e.g.*, SimCLR [5] and MoCo [17]), there has been an increasing interest in developing VLP methods to address a wide range of vision-and-language-related tasks. In general, VLP methods can be classified into two categories according to the vision-and-language interaction, *i.e.*, dual-encoders and fusion-encoders. Existing dual-encoder methods can be summarized according to the following aspects: (i) using medium-scale curated image-text data [48], (ii) using large-scale noisy image-text data [22], (iii) designing more fine-grained image-text contrast [63], (iv) adopting extra single modal contrastive learning [43]. For fusion-encoder approaches, existing studies can be further categorized with respect to these three perspectives: (i) Uni-modal encoders: different methods adopt different image features (*e.g.*, region features [28, 39], patch embeddings [25], and grid features [20]) and distinct text features (*e.g.*, statistic embeddings [25], and dynamic embeddings [16]); (ii) Multi-modal

---

[2]Although the terminologies "cross-modal" and "multi-modal" have been used interchangeably in the literature, we treat them as terms with different meanings in this paper.

[3]It is worth noting that most existing studies only conduct the evaluation on the vision-only tasks and disregard the language-only tasks although the text representations are simultaneously learned.

[4]Medical image-to-text generation refers to medical/radiology report generation in previous studies [10, 9].
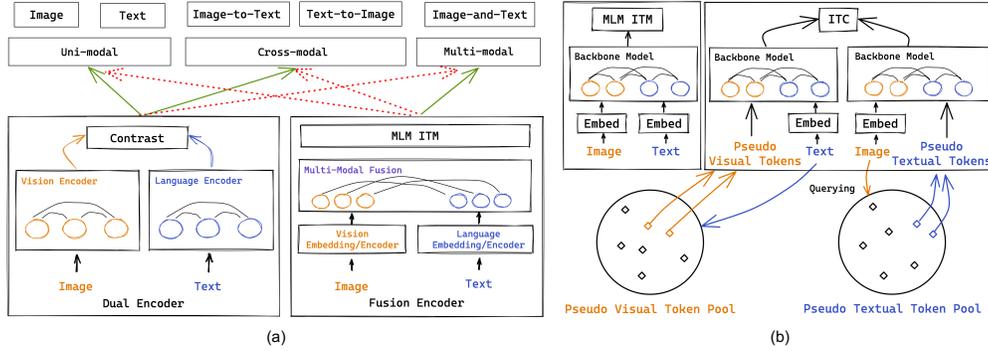
Figure 1: (a) Illustrations of two Med-VLP paradigms and their advantages (pointed by green arrows) and disadvantages (pointed by red arrows) in downstream tasks; (b) The overall architecture of our proposed approach, where the backbone models share the same parameters, and we duplicate them for illustration.

fusion modules: existing studies adopted the single-stream fusion scheme [53, 30] or dual-stream fusion scheme [55, 64]; (iii) Pretext tasks: existing studies explore a variety of pre-training tasks, including masked language modeling [28], masked image modeling [39, 6], image-text matching [65]. This paper adopts the model architecture of fusion encoders and the pretext tasks from both dual-encoder and fusion-encoder types.

**Medical Vision-and-Language Pre-Training (Med-VLP)**    Being one of the applications and extensions of VLP to the medical domain, Med-VLP aims to understand the content of medical images and texts, which can be traced back to [66] for dual-encoders and [31] for fusion-encoders. For dual-encoders, the follow-up studies [19, 44, 57] explored the global-local image-text contrastive learning to capture more fine-grained information among medical images and texts and have achieved state-of-the-art results in the medical image classification task. For fusion-encoders, [24, 42, 7] performed pre-training to improve the multi-modal reasoning ability of the vision-and-language models for the downstream task (*e.g.*, Medical VQA). Besides, [8] integrated medical knowledge into the pre-training procedure to improve the performance on downstream medical tasks.

**Unified Vision-and-Language Pre-training**    To unify the dual and fusion encoders, existing studies mainly adopted/designed specific model architectures to accommodate different pretext tasks. The most common scheme is to add an extra multi-modal fusion module to the dual encoders and perform the cross-modal pretext task (*i.e.*, image-text contrast) before the fusion and multi-modal pretext tasks (*e.g.*, MLM and ITM) after the fusion [27, 52]. Another line of research [11, 58] resorts to multi-tasking on various downstream supervised vision-language tasks by formulating them as sequence-to-sequence tasks. Besides, [4, 59] proposed a mixture-of-modality experts (MoME) Transformer to unify vision-and-language models by employing a set of modality experts to replace the feed-forward networks (FFN) in the standard Transformer. More recently, [14] proposed an encoder-decoder generative model learned from prefix language modeling and prefix image modeling. However, the aforementioned studies are either architecture-dependent or task-dependent, and they perform the unifying through training different parts of the models when applying different types of VLP objectives. Therefore, it is expected to unify the existing Med-VLP types in an *architecture- and task-agnostic* fashion to improve the generalization and extensionality ability of Med-VLP methods.

## 3   Bridging the Gap

In this section, we introduce the PTUnifier framework for unifying the fusion-encoder and dual-encoder types. §3.1 details the problem to be addressed. This work proposes to unify inputs using prompts (in §3.2). Thus one could jointly train various tasks even with different input formats (in §3.3) in either pre-training or fine-tuning.

### 3.1 Problem Definition

We adopt the general problem formulation for pre-training following existing studies [28, 55]. Formally, given a medical image $I$ and its corresponding description text $T$, the representation learning process can be formulated as

$$\theta^*, \theta_1^*, ..., \theta_S^* = \underset{\theta, \theta_1, ..., \theta_S}{\arg\min} \sum_{s=1}^{S} \mathcal{L}_s(Y_s, \mathcal{H}_{\theta_s}(\mathcal{M}_\theta(\boldsymbol{X}))), \tag{1}$$

where $S$ refers to the number of pretext tasks; $\mathcal{L}_s$ are the loss functions of pretext tasks; $Y_s$ are the corresponding ground-truth labels; $\mathcal{H}_{\theta_s}$ are the prediction heads with their parameters $\theta_s$; $\mathcal{M}_\theta$ is the backbone model which is parameterized by $\theta$; $\boldsymbol{X}$ represents the input to the backbone model, which could be one of the following cases:

$$\boldsymbol{X} = \begin{cases} (\boldsymbol{X}^v) & \text{if } \textit{image-only} \\ (\boldsymbol{X}^l) & \text{if } \textit{text-only} \\ (\boldsymbol{X}^v, \boldsymbol{X}^l) & \text{if } \textit{image-text} \end{cases} \tag{2}$$

where we suppose that we have embedded a medical image $I$ as $\boldsymbol{X}^v \in \mathbb{R}^{D_v \times N_v}$ or a medical text $T$ as $\boldsymbol{X}^l \in \mathbb{R}^{D_l \times N_l}$ when dealing with vision and language modalities. The challenge of the problem is to make the backbone model $\mathcal{M}_\theta$ deal with such variable-size and heterogeneous input. After overcoming this challenge, we can perform different types of downstream vision-language tasks (*i.e.*, uni-modal, cross-modal, and multi-modal tasks).

### 3.2 Unifying Inputs via Prompts

To unify inputs, we propose to unify the inputs via prompts so as to perform different types of tasks. In specific, we design two solutions, *i.e.*, a basic solution for *compatibility* and an advanced solution for *scalability*. In this work, we use the advanced solution in default if not specified.

#### 3.2.1 Compatibility using Prompts

To make the backbone model compatible with variable-size and heterogeneous input, this work proposes a simple yet effective approach, namely using Prompt (PT) as a placeholder for missing modality. $\mathcal{M}_\theta$ naturally accepts two inputs (visual and textual embeddings $\boldsymbol{X}^v, \boldsymbol{X}^l$), which is by definition compatible to inputs with image-text pairs. For image-only/text-only inputs, we propose to introduce visual/textual prompts to enable the backbone model to perceive the missing input in a specific modality:

$$\boldsymbol{X} = \begin{cases} (\boldsymbol{X}^v, \boldsymbol{PT}^l) & \text{if } \textit{image-only} \\ (\boldsymbol{PT}^v, \boldsymbol{X}^l) & \text{if } \textit{text-only} \\ (\boldsymbol{X}^v, \boldsymbol{X}^l) & \text{if } \textit{image-text} \end{cases} \tag{3}$$

where $\boldsymbol{PT}^v \in \mathbb{R}^{D_v \times k}$ and $\boldsymbol{PT}^l \in \mathbb{R}^{D_l \times k}$ are the visual and textual prompts, respectively.

#### 3.2.2 Scalability of Prompts

The above solution adopts a static fashion to introduce prompts, which might have limited diversity and therefore harm its capacity. Hence, we construct a pool of visual/textual prompts *instead of static prompts*. Importantly, the selection of prompts is *conditioned on the input embeddings*.

Formally, we define a visual prompt pool $\boldsymbol{V} \in \mathbb{R}^{D_v \times N_v}$ and a textual prompt pool $\boldsymbol{T} \in \mathbb{R}^{D_l \times N_l}$. $N_v$ and $N_l$ are the size of the visual/textual prompt pool, respectively. Given the image-only input with its visual embedding sequence $\boldsymbol{X}^v$ or language-only input with its textual embedding sequence $\boldsymbol{X}^l$, we conduct a pooling operation (*e.g.*, average/max pooling) to obtain a *query vector* for existing modality (denoted as $\boldsymbol{q}^v$ or $\boldsymbol{q}^l$), namely, $\boldsymbol{q}^v = \text{pooling}(\boldsymbol{X}^v)$ and $\boldsymbol{q}^l = \text{pooling}(\boldsymbol{X}^l)$, respectively. To get the prompts of the missing modality, the selection of prompts is based on the similarity scores between the query vector and all prompts in the pool from the missing modality:

$$\begin{aligned} \boldsymbol{PT}^l &= \underset{\boldsymbol{w} \in \boldsymbol{V}}{\text{top-}k} \left[ \boldsymbol{w}^T \boldsymbol{q}^v \right], \\ \boldsymbol{PT}^v &= \underset{\boldsymbol{w} \in \boldsymbol{T}}{\text{top-}k} \left[ \boldsymbol{w}^T \boldsymbol{q}^l \right], \end{aligned} \tag{4}$$

where $\boldsymbol{w}$ is an embedding vector in the prompt pool, and we select $k$ closest prompts as the input embedding sequence of the missing modality.

**Intuitive Explanation**    Without loss of any generality, we take a text-only scenario as an example, but it also holds for the image-only scenario. To select the best visual prompts for the text-only input, the proposed method chooses the most similar ones compared to the given textual query vector. As an intuitive explanation, one could treat the visual prompt pool as a feature bank that stores the most representative images of a given dataset, Eq. 4 aims to choose the visual prompts that might convey a similar semantic meaning as the given text by conducting dot products. In other words, *it might, at least to some extent, automatically fill (originally unprovided) semantically-similar images conditioned on purely the given text.*

**Linking to Prompts**    We find that the PTUnifier (especially the static one in §3.2.1) is quite similar to the prompt tuning [29, 35]. They both introduce special tokens or vectors as a certain signal for training or inference. One notable difference is that in a special version of PTUnifier using prompts pools (see §3.2.2), the selection of additional tokens/vectors is conditioned on the input, while prompts are generally static and constant to input.

### 3.3    Unifying Multiple Pre-training Objectives

Owing to the unified image and/or text input formulation, we can adopt pretext tasks of both fusion-encoders and dual-encoders (see Eq. 1). Following previous studies [28, 55, 66, 48], we develop two commonly used pretext tasks (*i.e.*, masked language modeling (MLM) and image-text matching (ITM)) for fusion-encoders and the image-text contrast (ITC) pretext task for dual-encoders. To produce the prediction for the aforementioned MLM and ITM tasks, we use two independent prediction heads $\mathcal{H}_{\mathrm{MLM}}$ and $\mathcal{H}_{\mathrm{ITM}}$ (*i.e.*, two two-layer multilayer perceptrons (MLP)).

**Masked Language Modeling (MLM)**    Following BERT [13], we randomly mask 15% of the words (denoted as $Y_{\mathrm{MLM}}$) of the input text $T$ and recover them according to the remaining text ($T_{\mathrm{M}}$) and the input $I$. The MLM objective is given by:

$$\mathcal{L}_{\mathrm{MLM}} = - \sum_{(I,T)} \log p_{\mathrm{MLM}}(Y_{\mathrm{MLM}}|I, T_{\mathrm{M}}), \qquad (5)$$

where $p_{\mathrm{MLM}}$ is obtained by applying $\mathcal{H}_{\mathrm{MLM}}$ followed by a softmax operation on the corresponding representations of [MASK] in $\boldsymbol{Z}^{l}$.

**Image-Text Matching (ITM)**    aims to distinguish whether an image-text pair is a match. In detail, a positive image-text pair and a randomly sampled negative pair are fed into $\mathcal{M}_{\theta}$ and the concatenation of $\boldsymbol{z}^{v}_{[\mathrm{CLS}]}$ and $\boldsymbol{z}^{l}_{[\mathrm{CLS}]}$ is processed by $\mathcal{H}_{\mathrm{MLM}}$ followed by a softmax layer to output a binary probability $p_{\mathrm{ITM}}$. Therefore, the ITM objective is given by

$$\mathcal{L}_{\mathrm{ITM}} = - \sum_{(I,T)} \log p_{\mathrm{ITM}}(Y_{\mathrm{ITM}}|I, T). \qquad (6)$$

**Image-Text Contrast (ITC)**    aims to learn better uni-modal/cross-modal representation from the instance-level contrast. In this work, given an image-text pair, we use two different forward procedures on the image-only input $I$ and the text-only input $T$, respectively, to obtain the image-only representation (denoted as $\boldsymbol{z}^{v}$) and text-only representation (denoted as $\boldsymbol{z}^{l}$). Afterward, we adopt the similarity function $s(I, T) = \boldsymbol{z}^{v\top}\boldsymbol{z}^{l}$ to compute the image-to-text similarity and text-to-image similarity between $\boldsymbol{z}^{v}$ and $\boldsymbol{z}^{l}$. Subsequently, the similarities are normalized as follows:

$$p^{\mathrm{i2t}}_{n} = \frac{\exp\left(s\left(I, T_{n}\right)/\tau\right)}{\sum_{n=1}^{N} \exp\left(s\left(I, T_{n}\right)/\tau\right)}, \qquad (7)$$

$$p^{\mathrm{t2i}}_{n} = \frac{\exp\left(s\left(I_{n}, T\right)/\tau\right)}{\sum_{n=1}^{N} \exp\left(s\left(I_{n}, T\right)/\tau\right)}, \qquad (8)$$

where $N$ is the size of the mini-batch. The ground-truth labels $Y^{\mathrm{i2t}}$ and $Y^{\mathrm{t2i}}$ are two $N \times N$ one-hot matrices, where negative pairs have a probability of 0 and the positive pair has a probability of 1.

Therefore, the ITC objective is given by

$$\mathcal{L}_{ITC} = -\frac{1}{2} \sum_{(I,T)} \log p^{\text{i2t}}(Y^{\text{i2t}}|I,T) - \frac{1}{2} \sum_{(I,T)} \log p^{\text{t2i}}(Y^{\text{t2i}}|I,T). \tag{9}$$

## 4   The Model Architecture

The previous section documents the unification at the input and task levels. This section will introduce the overall architecture of our work. As a pipeline, we first map visual and textual tokens into embeddings space ($\boldsymbol{X}^v$ and $\boldsymbol{X}^l$ as specified in §4.1). Such token embeddings with or without prompts will be jointly processed by an identical backbone model $\mathcal{M}_\theta$ (§4.2). An overview of the proposed approach is shown in Figure 1(b).

### 4.1   Visual and Textual Embeddings

**Visual embedding**   For an input image $I$, it is first segmented into patches following [15]. Then the patches are linearly projected into patch embeddings $\boldsymbol{X}^v = (\boldsymbol{x}_1^v, \boldsymbol{x}_2^v, \ldots, \boldsymbol{x}_{N_v}^v), \boldsymbol{x}_i^v \in \mathbb{R}^{D_v}$ through a linear transformation and a special learnable token embedding $\boldsymbol{x}_{[\text{CLS}]}^v$ is prepended for the aggregation of visual information. Therefore, the image embedding sequence is obtained by summing up the patch embeddings and learnable 1D position embeddings $\boldsymbol{E}_{pos}^v \in \mathbb{R}^{D_v \times (N_v+1)}$:

$$\boldsymbol{X}^v = [\boldsymbol{x}_{[\text{CLS}]}^v; \boldsymbol{x}_1^v; \boldsymbol{x}_2^v; ...; \boldsymbol{x}_{N_v}^v] + \boldsymbol{E}_{pos}^v, \tag{10}$$

where $[\cdot;\cdot]$ represents the column concatenation.[5]

**Textual embedding**   Similarly, for an input text $T$, we follow BERT [13] to tokenize the input text to subword tokens by WordPiece [62]. Afterwards, the tokens are linearly projected into embeddings $\boldsymbol{X}^l = (\boldsymbol{x}_1^l, \boldsymbol{x}_2^l, ..., \boldsymbol{x}_{N_l}^l), \boldsymbol{x}_i^l \in \mathbb{R}^D$ through a linear transformation with a start-of-sequence token embedding $\boldsymbol{x}_{[\text{CLS}]}^l$, and a special boundary token embedding $\boldsymbol{x}_{[\text{SEP}]}^l$ added. Therefore, the text embedding sequence is obtained by summing up the sub-word token embeddings and text position embeddings $\boldsymbol{E}_{pos}^l \in \mathbb{R}^{D \times (N_l+2)}$:

$$\boldsymbol{X}^l = [\boldsymbol{x}_{[\text{CLS}]}^l; \boldsymbol{x}_1^l; \ldots; \boldsymbol{x}_{N_l}^l; \boldsymbol{x}_{[\text{SEP}]}^l] + \boldsymbol{E}_{pos}^l. \tag{11}$$

### 4.2   The Backbone Model

Since the input image and/or text are represented as a unified image-text sequence, the backbone model can be any model for sequential modeling. In this work, we adopt an attention-based Med-VLP model with the multi-modal interaction, which can be an effective model (including uni-modal encoders and a multi-modal fusion module) or an efficient one (*i.e.*, a single Transformer model), where the attention mechanism is defined as

$$\text{ATTN}(\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}) = \text{softmax}\left(\boldsymbol{Q}\boldsymbol{K}^T/\sqrt{D_k}\right)\boldsymbol{V}, \tag{12}$$

where $\boldsymbol{Q}$, $\boldsymbol{K}$, and $\boldsymbol{V}$ are the query, key, and value matrix linearly transformed from the input embedding sequence, respectively, and $D_k$ is the dimension of $\boldsymbol{K}$. Formally, for a given input (defined in Eq. 3), the whole representation process can be formulated as

$$\boldsymbol{Z}^v, \boldsymbol{Z}^l = \mathcal{M}_\theta(\boldsymbol{X}), \tag{13}$$

where $\boldsymbol{Z}^v = (\boldsymbol{z}_{[\text{CLS}]}^v, \boldsymbol{z}_1^v, \boldsymbol{z}_2^v, ..., \boldsymbol{z}_{N_v}^v)$ and $\boldsymbol{Z}^l = (\boldsymbol{z}_{[\text{CLS}]}^l, \boldsymbol{z}_1^l, ..., \boldsymbol{z}_{N_l}^l, \boldsymbol{z}_{[\text{SEP}]}^l)$ are the image and text representations from the backbone model.

## 5   Experimental Settings

### 5.1   Pre-training Datasets

In our experiments, we perform the pre-training on three datasets, which are described as follows:

---

[5]We overload the notation $\boldsymbol{X}^v$ for simplicity (same for $\boldsymbol{X}^l$).

- **ROCO** [47]: a dataset of radiology figure-caption pairs from PubMed Central, an open-access biomedical literature database.
- **MedICaT** [54]: a dataset of medical figure-caption pairs also extracted from PubMed Central. Different from ROCO, 75% of its figures are compound figures, including several sub-figures.
- **MIMIC-CXR** [23]: the largest radiology dataset to date consisting of chest X-ray images (in frontal or lateral views) and their reports from the Beth Israel Deaconess Medical Center.

For all the datasets, we exclude those samples with a text length of less than 3. For ROCO and MedICaT, we filter non-radiology samples, and for MIMIC-CXR, we only keep images in the frontal view. As for the dataset split, we adopt the official splits of ROCO and MIMIC-CXR. For MedICaT, we randomly sample 1,000 image-text pairs for validation and 1,000 for testing, and the remaining image-text pairs are used for training. Different from the texts in general-domain VLP, medical texts are long narratives consisting of multiple sentences. To deal with this case, we randomly sample a sentence from the input text in each iteration.

## 5.2 Medical Vision-Language Benchmark

To evaluate the performance, we construct a medical vision-language evaluation benchmark including three types of tasks, *i.e.*, uni-modal, cross-modal, and multi-modal evaluations.[6] All the adopted datasets are related to radiology.

**Uni-modal Evaluation** requires the model to process a single modality with vision-only or language-only inputs. For vision-only tasks, we conduct the image classification (IC) experiments on CheXpert [21] and RSNA Pneumonia [51]. For language-only tasks, we perform both the understanding task (*i.e.*, text classification (TC)) and the generation task (*i.e.*, text summarization (TS)) on the RadNLI [49, 41] and MIMIC-CXR datasets, respectively.

**Cross-modal Evaluation** requires the model to align the vision and language modalities. We conduct experiments on three kinds of tasks (*i.e.*, image-to-text retrieval (ITR), text-to-image retrieval (TIR), and image-to-text generation (ITG)). For ITR and TIR, we adopt the ROCO dataset and measure both zero-shot and fine-tuned performance. During the evaluation, we sample 2,000 image-text pairs from the ROCO test set and report the results on the 2,000 sampled image-text pairs due to the large time complexity of the ranking process. For ITG, we conduct experiments on the MIMIC-CXR dataset to evaluate its ability for radiology report generation.

**Multi-modal Evaluation** requires the model to reason over both the image and text inputs through the multi-modal interaction. We conduct the experiments on the medical visual question answering (VQA) task, which requires the model to answer natural language questions about a medical image. We adopt three publicly available Med-VQA datasets (*i.e.*, VQA-RAD [26], SLAKE [34], and MedVQA-2019 [1]), where VQA-RAD consists of 3,515 image-question pairs, SLAKE contains 14,028 image-question pairs and MedVQA-2019 contains 15,292 image-question pairs.

## 5.3 Implementation Details

**Pre-training** We adopt the classical VLP model as the backbone model, including a vision encoder, a language encoder, and a multi-modal fusion module. For the vision and language encoders, we adopt base-size Transformer encoders with 12 layers initialized from CLIP-ViT-B [48] RoBERTa-base [37] and their hidden dimension is set to 768. For the multi-modal fusion module, we set the number of Transformer layers to 6, the dimension of the hidden states to 768, and the number of heads to 12. For the visual/textual prompt pools, the dimension and the pool size is set to 768 and 1,024, respectively, by default. For optimization, the pre-training takes 100,000 steps with AdamW optimizer [38] with a weight decay of 0.01. The learning rates for the vision and language encoders and the remaining parameters are set to 1e-5 and 5e-5, respectively. We use the warm-up strategy during the first 10% of the total number of steps, and the learning rate is linearly decayed to 0 after warm-up. For data augmentation, we use center-crop to resize each image to the size of $288 \times 288$.

**Fine-tuning** For all downstream tasks, we use the AdamW optimizer with the learning rate set to 5e-6 and 2.5e-4 for the backbone model and task-specific layers, respectively. The fine-tuning

---

[6]More details of the downstream evaluations are reported in Appendix A.

| Methods | Uni-Modal | | | | Cross-Modal | | | Multi-Modal | | |
| | Image | | Text | | Image-to-Text | | Text-to-Image | | | |
| | CheXpert AUROC | PNAS AUROC | RadNLI Acc | MIMIC RL | MIMIC BL4 | ROCO R@1 | ROCO R@1 | VQA-RAD Acc | SLAKE Acc | MedVQA-2019 Acc |
| Study$_1$ | ConVIRT [66] | | ClinicalBERT [2] | TransABS [36] | R2Gen [10] | | ViLT [25] | | CPRD [33] | |
| | 87.3 | 81.3 | 72.6 | 43.8 | 8.0 | 11.9 | 9.8 | 72.7 | 82.1 | - |
| Study$_2$ | GLoRIA [19] | | IFCC [41] | WGSum [18] | M2Trans [41] | | METER [16] | | MMBERT [24] | |
| | 88.1 | 88.6 | 77.8 | 45.1 | 10.5 | 14.5 | 11.3 | 72.0 | - | 77.9 |
| PTUnifier (ours) | 90.1 | 90.6 | 80.0 | 46.2 | 10.7 | 21.0 | 20.8 | 78.3 | 85.2 | 79.3 |

Table 1: Comparisons of our proposed method with previous studies on three types of evaluations (*i.e.*, uni-modal, cross-modal, and multi-modal evaluations). Study$_1$ and Study$_2$ denote two state-of-the-art approaches of each type of tasks, respectively. BL-4 denotes BLEU score using 4-grams and RG-L denotes ROUGE-L (same below). Dark and light grey colors highlight the top and second best results on each metric (same below). Note that the results of text summarization and image-to-text generation are replicated using our pre-processed data (See Appendix A).

strategies can be divided into three categories according to the type of tasks. Specifically, for the classification tasks (*i.e.*, IC, TC, and VQA), we feed the concatenation of the image/visual prompt and text/textual prompt representations to a randomly initialized two-layer MLP to predict the labels. For the retrieval tasks (*i.e.*, ITR and TIR), we adopt the prediction head for the image-text contrast pre-text task and test its zero-shot and fine-tuned performance. For the generation tasks (*i.e.*, TS and ITG), we feed the concatenation of the sequence of image/visual prompt and text/textual prompt representations to a Transformer decoder with its parameters (except for the parameters of cross-attention layers) initialized from the pre-trained language encoder. For the evaluation metrics, we follow the previous studies to adopt AUROC for IC, accuracy for TC and VQA, Recall@K (K=1, 5, 10) for ITR and TIR, and natural language generation (NLG) metrics (*i.e.*, BLEU [46], METEOR [12], CIDEr [56], and ROUGE [32]) for TS and ITG.

All pre-training and fine-tuning experiments are conducted on 80GB NVIDIA A100 GPUs with mixed-precision [40] to accelerate training and save memory.

## 6 Results and Analyses

### 6.1 Main Results

To demonstrate the effectiveness of the proposed approach, we conduct experiments on the aforementioned medical vision-language benchmark. The results of the main experiments are reported in Table 1. There are several observations. First, our approach achieves the best performance on all tasks. It outperforms previous studies on uni-modal image classification (+2.0% AUROC), text classification (+2.2% Accuracy), text summarization (+1.1% Rouge-L), image-to-text generation (+0.2% BLEU-4), image-to-text retrieval (+7.5% Recall@1), text-to-image retrieval (+9.5% Recall@1), and multi-modal VQA (+3.4% Accuracy), which confirms the validity of the proposed approach. Second, the proposed approach outperforms those complicated methods designed for specific tasks. For example, R2Gen introduced recurrent memory networks to the Transformer decoder to augment its decoding ability; WGSum used extra word graphs to improve the ability to detect keywords in the findings section of radiology reports. CPRD adopted representation distillation to alleviate the data scarcity problem in the Med-VQA task. These observations show that different pre-training ways can enable distinct abilities of the model, and it is possible to design an appropriate approach to exploit the knowledge shared across different tasks and perform various tasks using a unified model. Note that *the existing studies are only designed for a single task*, while our approach generally targets all vision- and/or language-related tasks, namely, without any tailored adaptations to a specific task.

### 6.2 Ablation Study

To further illustrate the effectiveness of our proposed approach, we perform an ablation study on the pre-training objectives, including the ones from fusion-encoders (*i.e.*, MLM and ITM) and the one from dual-encoders (*i.e.*, ITC).

| | Objectives | | | Uni-Modal Image | | | Text | Cross-Modal Image-to-Text | | Multi-Modal | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ID | | | | | CheXpert | | RadNLI | MIMIC | | | VQA-RAD | |
| | MLM | ITM | ITC | 1% AUROC | 10% AUROC | 100% AUROC | Acc | BL4 | CDr | Open Acc | Closed Acc | Overall Acc |
| 1 | ✓ | | | 66.1 | 79.1 | 81.1 | 77.2 | 6.9 | 11.1 | 57.5 | 79.5 | 70.8 |
| 2 | | ✓ | | 56.9 | 83.0 | 85.8 | 77.5 | 10.0 | 18.2 | 23.5 | 82.8 | 59.3 |
| 3 | ✓ | ✓ | | 74.5 | 87.2 | 88.4 | 78.3 | 9.9 | 17.1 | 67.0 | 84.6 | 77.7 |
| 4 | | | ✓ | 88.0 | 88.9 | 89.3 | 76.5 | 10.3 | 19.0 | 64.8 | 81.0 | 74.6 |
| 5 | ✓ | ✓ | ✓ | 88.7 | 89.0 | 90.1 | 80.0 | 10.7 | 21.0 | 68.7 | 84.6 | 78.3 |

Table 2: Ablation studies on the different types of objectives, including the fusion-encoders ones (*i.e.*, masked language modeling (MLM) and image-text matching (ITM)) and the dual-encoders one (*i.e.*, image-text contrast (ITC)). 1%, 10%, and 100% represent the different portion of training data.

There are several observations drawn from different aspects. First, the objectives of fusion encoders (*i.e.*, MLM and ITM) guide the models (*i.e.*, ID 3 and 5) to the more powerful multi-modal representations than other models without them, which could be observed from the performance on the downstream Med-VQA task. Second, the image-text contrast objective of dual encoders assists the models (*i.e.*, ID 4 and 5) in learning the uni-modal image representations and the cross-modal representations, and the models pre-trained with the ITC objective outperform those pre-trained without the ITC objective. More importantly, the models pre-trained with the ITC objective (*i.e.*, ID 4 and 5) demonstrate their great transfer ability where the pre-trained models can achieve high performance with very little data (*e.g.*, 1% and 10%). Third, it is interesting to note that the ITC objective does not promote the performance of the uni-modal text classification task. We can explain this phenomenon by the reason that images and texts are abstracted at different levels, where pixels of images have a lower semantic level than tokens of texts. Therefore, in the ITC process, the texts can be treated as a kind of "supervision signals" for the learning of image encoding, yet, it is harder for the images to play such a role in contrast. This can be observed from previous studies [48, 22, 43], where the dual-encoders were only evaluated on the uni-modal vision tasks or cross-modal tasks. Fourth, performing both types of objectives promotes the model (*i.e.*, ID 5) to achieve the best performance across all the tasks, which confirms the feasibility of the research direction on unifying the fusion-encoders and dual-encoders.

## 7 Conclusion

In this paper, we proposed a simple yet effective scheme to take advantage of both fusion encoders, and dual encoders, where visual and textual prompt pools are used to make our model compatible with different kinds of inputs (*i.e.*, image-only, text-only, and image-text-pair), and thus different types of objectives (*e.g.*, MLM and ITM for fusion-encoders and ITC for dual-encoders) can be adopted for pre-training. It is worth noting that our proposed approach is complementary to most of the existing Med-VLP models. Experimental results confirm the validity of our approach, where it achieves state-of-the-art performance on the downstream tasks. The further analyses investigate the effects of different types of objectives. Such empirical studies might provide a valuable reference for future research in this area.

## References

[1] Asma Ben Abacha, Sadid A. Hasan, Vivek Datla, Joey Liu, Dina Demner-Fushman, and Henning Müller. Vqa-med: Overview of the medical visual question answering task at imageclef 2019. In *CLEF*, 2019. 7

[2] Emily Alsentzer, John Murphy, William Boag, Wei-Hung Weng, Di Jindi, Tristan Naumann, and Matthew McDermott. Publicly available clinical bert embeddings. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 72–78, 2019. 8

[3] Shruthi Bannur, Stephanie Hyland, Qianchu Liu, Fernando Perez-Garcia, Maximilian Ilse, Daniel C Castro, Benedikt Boecking, Harshita Sharma, Kenza Bouzid, Anja Thieme, et al. Learning to exploit temporal structure for biomedical vision-language processing. *arXiv preprint arXiv:2301.04558*, 2023. 2

[4] Hangbo Bao, Wenhui Wang, Li Dong, Qiang Liu, Owais Khan Mohammed, Kriti Aggarwal, Subhojit Som, Songhao Piao, and Furu Wei. Vlmo: Unified vision-language pre-training with mixture-of-modality-experts. In *Advances in Neural Information Processing Systems*, 2022. 2, 3

[5] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 2

[6] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *European conference on computer vision*, pages 104–120. Springer, 2020. 3

[7] Zhihong Chen, Yuhao Du, Jinpeng Hu, Yang Liu, Guanbin Li, Xiang Wan, and Tsung-Hui Chang. Multi-modal masked autoencoders for medical vision-and-language pre-training. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part V*, pages 679–689. Springer, 2022. 2, 3

[8] Zhihong Chen, Guanbin Li, and Xiang Wan. Align, reason and learn: Enhancing medical vision-and-language pre-training with knowledge. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 5152–5161, 2022. 2, 3

[9] Zhihong Chen, Yaling Shen, Yan Song, and Xiang Wan. Cross-modal memory networks for radiology report generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5904–5914, 2021. 2

[10] Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. Generating radiology reports via memory-driven transformer. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, pages 1439–1449, 2020. 2, 8

[11] Jaemin Cho, Jie Lei, Hao Tan, and Mohit Bansal. Unifying vision-and-language tasks via text generation. In *International Conference on Machine Learning*, pages 1931–1942. PMLR, 2021. 3

[12] Michael Denkowski and Alon Lavie. Meteor 1.3: Automatic metric for reliable optimization and evaluation of machine translation systems. In *Proceedings of the sixth workshop on statistical machine translation*, pages 85–91, 2011. 8

[13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, 2019. 2, 5, 6

[14] Shizhe Diao, Wangchunshu Zhou, Xinsong Zhang, and Jiawei Wang. Prefix language models are unified modal learners. *arXiv preprint arXiv:2206.07699*, 2022. 3

[15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020. 6

[16] Zi-Yi Dou, Yichong Xu, Zhe Gan, Jianfeng Wang, Shuohang Wang, Lijuan Wang, Chenguang Zhu, Zicheng Liu, Michael Zeng, et al. An empirical study of training end-to-end vision-and-language transformers. *arXiv preprint arXiv:2111.02387*, 2021. 2, 8

[17] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 2

[18] Jinpeng Hu, Jianling Li, Zhihong Chen, Yaling Shen, Yan Song, Xiang Wan, and Tsung-Hui Chang. Word graph guided summarization for radiology findings. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4980–4990, 2021. 8

[19] Shih-Cheng Huang, Liyue Shen, Matthew P Lungren, and Serena 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*, pages 3942–3951, 2021. 2, 3, 8, 14

[20] Zhicheng Huang, Zhaoyang Zeng, Bei Liu, Dongmei Fu, and Jianlong Fu. Pixel-bert: Aligning image pixels with text by deep multi-modal transformers. *arXiv preprint arXiv:2004.00849*, 2020. 2

[21] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, et al. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 590–597, 2019. 7

[22] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pages 4904–4916. PMLR, 2021. 2, 9

[23] Alistair EW Johnson, Tom J Pollard, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Yifan Peng, Zhiyong Lu, Roger G Mark, Seth J Berkowitz, and Steven Horng. Mimic-cxr-jpg, a large publicly available database of labeled chest radiographs. *arXiv preprint arXiv:1901.07042*, 2019. 2, 7

[24] Yash Khare, Viraj Bagal, Minesh Mathew, Adithi Devi, U Deva Priyakumar, and CV Jawahar. Mmbert: Multimodal bert pretraining for improved medical vqa. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pages 1033–1036. IEEE, 2021. 2, 3, 8

[25] Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convolution or region supervision. In *International Conference on Machine Learning*, pages 5583–5594. PMLR, 2021. 2, 8

[26] Jason J Lau, Soumya Gayen, Asma Ben Abacha, and Dina Demner-Fushman. A dataset of clinically generated visual questions and answers about radiology images. *Scientific data*, 5(1):1–10, 2018. 7

[27] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in Neural Information Processing Systems*, 34, 2021. 3

[28] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019. 2, 3, 4, 5

[29] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, 2021. 5

[30] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *European Conference on Computer Vision*, pages 121–137. Springer, 2020. 3

[31] Yikuan Li, Hanyin Wang, and Yuan Luo. A comparison of pre-trained vision-and-language models for multimodal representation learning across medical images and reports. In *2020 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 1999–2004. IEEE, 2020. 2, 3

[32] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004. 8

[33] Bo Liu, Li-Ming Zhan, and Xiao-Ming Wu. Contrastive pre-training and representation distillation for medical visual question answering based on radiology images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 210–220. Springer, 2021. 8

[34] Bo Liu, Li-Ming Zhan, Li Xu, Lin Ma, Yan Yang, and Xiao-Ming Wu. Slake: A semantically-labeled knowledge-enhanced dataset for medical visual question answering. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pages 1650–1654. IEEE, 2021. 7, 14

[35] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35, 2023. 5

[36] Yang Liu and Mirella Lapata. Text summarization with pretrained encoders. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3730–3740, 2019. 8

[37] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019. 7

[38] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. 7

[39] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems*, 32, 2019. 2, 3

[40] Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. Mixed precision training. In *International Conference on Learning Representations*, 2017. 8

[41] Yasuhide Miura, Yuhao Zhang, Emily Tsai, Curtis Langlotz, and Dan Jurafsky. Improving factual completeness and consistency of image-to-text radiology report generation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5288–5304, 2021. 7, 8, 14

[42] Jong Hak Moon, Hyungyung Lee, Woncheol Shin, Young-Hak Kim, and Edward Choi. Multi-modal understanding and generation for medical images and text via vision-language pre-training. *IEEE Journal of Biomedical and Health Informatics*, 26(12):6070–6080, 2022. 2, 3

[43] Norman Mu, Alexander Kirillov, David Wagner, and Saining Xie. Slip: Self-supervision meets language-image pre-training. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVI*, pages 529–544. Springer, 2022. 2, 9

[44] Philip Müller, Georgios Kaissis, Congyu Zou, and Daniel Rückert. Joint learning of localized representations from medical images and reports. *ArXiv preprint*, abs/2112.02889, 2021. 2, 3

[45] Binh D Nguyen, Thanh-Toan Do, Binh X Nguyen, Tuong Do, Erman Tjiputra, and Quang D Tran. Overcoming data limitation in medical visual question answering. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 522–530. Springer, 2019. 14

[46] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002. 8

[47] Obioma Pelka, Sven Koitka, Johannes Rückert, Felix Nensa, and Christoph M Friedrich. Radiology objects in context (roco): a multimodal image dataset. In *Intravascular Imaging and Computer Assisted Stenting and Large-Scale Annotation of Biomedical Data and Expert Label Synthesis*, pages 180–189. Springer, 2018. 2, 7

[48] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 2, 5, 7, 9

[49] Alexey Romanov and Chaitanya Shivade. Lessons from natural language inference in the clinical domain. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1586–1596, 2018. 7

[50] Dhruv Sharma, Sanjay Purushotham, and Chandan K Reddy. Medfusenet: An attention-based multimodal deep learning model for visual question answering in the medical domain. *Scientific Reports*, 11(1):1–18, 2021. 14

[51] George Shih, Carol C Wu, Safwan S Halabi, Marc D Kohli, Luciano M Prevedello, Tessa S Cook, Arjun Sharma, Judith K Amorosa, Veronica Arteaga, Maya Galperin-Aizenberg, et al. Augmenting the national institutes of health chest radiograph dataset with expert annotations of possible pneumonia. *Radiology: Artificial Intelligence*, 1(1):e180041, 2019. 7

[52] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. Flava: A foundational language and vision alignment model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15638–15650, 2022. 2, 3

[53] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. Vl-bert: Pre-training of generic visual-linguistic representations. In *International Conference on Learning Representations*, 2019. 3

[54] Sanjay Subramanian, Lucy Lu Wang, Ben Bogin, Sachin Mehta, Madeleine van Zuylen, Sravanthi Parasa, Sameer Singh, Matt Gardner, and Hannaneh Hajishirzi. Medicat: A dataset of medical images, captions, and textual references. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2112–2120, 2020. 2, 7

[55] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5100–5111, 2019. 3, 4, 5

[56] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575, 2015. 8

[57] Fuying Wang, Yuyin Zhou, Shujun Wang, Varut Vardhanabhuti, and Lequan Yu. Multi-granularity cross-modal alignment for generalized medical visual representation learning. In *Advances in Neural Information Processing Systems*, 2022. 2, 3

[58] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning*, pages 23318–23340. PMLR, 2022. 3

[59] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, et al. Image as a foreign language: Beit pretraining for all vision and vision-language tasks. *arXiv preprint arXiv:2208.10442*, 2022. 3

[60] Zifeng Wang, Zhenbang Wu, Dinesh Agarwal, and Jimeng Sun. Medclip: Contrastive learning from unpaired medical images and text. *arXiv preprint arXiv:2210.10163*, 2022. 2

[61] Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Medklip: Medical knowledge enhanced language-image pre-training. *medRxiv*, pages 2023–01, 2023. 2

[62] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*, 2016. 6

[63] Lewei Yao, Runhui Huang, Lu Hou, Guansong Lu, Minzhe Niu, Hang Xu, Xiaodan Liang, Zhenguo Li, Xin Jiang, and Chunjing Xu. Filip: Fine-grained interactive language-image pre-training. In *International Conference on Learning Representations*, 2022. 2

[64] Fei Yu, Jiji Tang, Weichong Yin, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. Ernie-vil: Knowledge enhanced vision-language representations through scene graphs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 3208–3216, 2021. 3

[65] Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. Vinvl: Revisiting visual representations in vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5579–5588, 2021. 3

[66] Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D Manning, and Curtis P Langlotz. Contrastive learning of medical visual representations from paired images and text. *arXiv preprint arXiv:2010.00747*, 2020. 2, 3, 5, 8

# A  More Details of Downstream Evaluation

In this section, we detail the descriptions for each downstream evaluation dataset.

**CheXpert**    This dataset contains 224,316 chest radiographs labeled for 14 medical observations. Following the previous studies (GLoRIA [19]), we only keep those front-view radiographs, hold out the expert-labeled validation set as the test set, and randomly sample 5,000 images from the training data for validation.

**RNAS Pneumonia**    This dataset consists of 30,000 front-view chest radiographs labeled by "pneumothorax negative" or "pneumothorax positive". Following the previous studies (GLoRIA [19]), the train/validation/test split constitutes 70%/15%/15% of the dataset, respectively.

**RadNLI**    This dataset contains 19k sentence pairs labeled by "Entailment", "Neutral", or "Contradiction". We follow IFCC [41] to produce and pre-process the dataset, which contains the training data from an extra NLI dataset (*i.e.*, MedNLI).

**ROCO**    This dataset contains 81k image-text pairs. For the training and validation set, we adopt the official ones. For the test procedure, we sample 2,000 pairs from the test set and evaluate the models on the 2,000 pairs to obtain the Recall@K scores.

**MIMIC-CXR**    This dataset contains 377,110 chest X-rays. Different from the pre-training, for downstream evaluations (*i.e.*, text summarization and image-to-text generation), we only keep those front-view x-rays with both the findings and impression section.

**VQA-RAD**    This dataset consists of 315 images and 3,515 questions. We adopt the commonly used version pre-processed by MEVF [45].

**SLAKE**    This dataset contains 642 images and 14,028 questions. We follow the original SLAKE paper [34] to prepare and pre-process the dataset and adopt the official dataset split.

**MedVQA-2019**    This dataset contains 4,200 images and 15,292 questions. We follow previous studies [50] to prepare and pre-process the dataset by keeping the main three categories of questions: Modality, Plane, and Organ system.