

PMC-CLIP: Contrastive Language-Image Pre-training using Biomedical Documents

Weixiong Lin^{1,*}, Ziheng Zhao^{1,*}, Xiaoman Zhang^{1,2}, Chaoyi Wu^{1,2}, Ya Zhang^{1,2}, Yanfeng Wang^{1,2}, and Weidi Xie^{1,2,†}

¹ Cooperative Medianet Innovation Center, Shanghai Jiao Tong University, Shanghai, China

² Shanghai AI Laboratory, Shanghai, China

{wx.lin,Zhao.Ziheng,xm99sjtu,wtzxxxwcy02,ya.zhang,wangyanfeng,weidi}@sjtu.edu.cn

Abstract. Foundation models trained on large-scale dataset gain a recent surge in CV and NLP. In contrast, development in biomedical domain lags far behind due to data scarcity. To address this issue, we build and release PMC-OA, a biomedical dataset with 1.6M image-caption pairs collected from PubMedCentral’s OpenAccess subset, which is 8 times larger than before. PMC-OA covers diverse modalities or diseases, with majority of the image-caption samples aligned at finer-grained level, *i.e.*, subfigure and subcaption. While pretraining a CLIP-style model on PMC-OA, our model named PMC-CLIP achieves state-of-the-art results on various downstream tasks, including image-text retrieval on ROCO, MedMNIST image classification, Medical VQA, *i.e.* +8.1% R@10 on image-text retrieval, +3.9% accuracy on image classification.

Keywords: Foundation Model · Multimodal Dataset · Vision-Language Pretraining.

1 Introduction

In the recent literature, development of foundational models has been the main driving force in artificial intelligence, for example, large language models [26,8,2,23] trained with either autoregressive prediction or masked token inpainting, and computer vision models [17,25,30] trained by contrasting visual-language features. In contrast, development in the biomedical domain lags far behind due to limitations of data availability from two aspects, (i) the expertise required for annotation, (ii) privacy concerns. This paper presents our preliminary study for constructing a **large-scale, high-quality, image-text** biomedical dataset using publicly available scientific papers, with **minimal manual efforts** involved.

In particular, we crawl figures and corresponding captions from scientific documents on PubMed Central, which is a free full-text archive of biomedical and life sciences journal literature at the U.S. National Institutes of Health’s

†: Corresponding author. *: These authors contribute equally to this work.

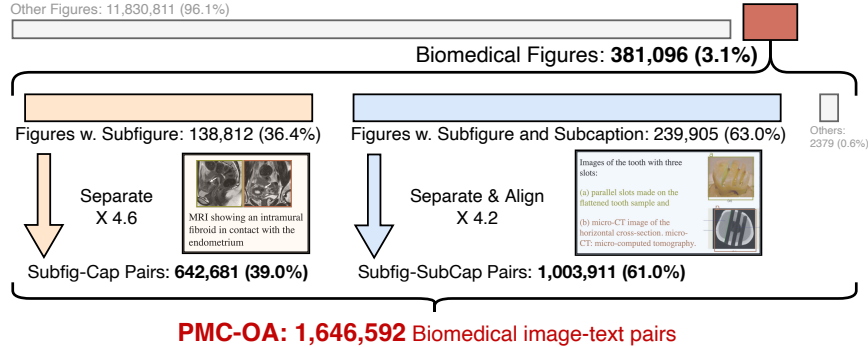


Fig. 1. Statistics over the pipeline and the collected PMC-OA.

National Library of Medicine (NIH/NLM) [27]. This brings two benefits: (i) the contents in publications are generally well-annotated and examined by experts, (ii) the figures have been well-anonymized and de-identified. In the literature, we are clearly not the first to construct biomedical datasets in such a manner, however, existing datasets suffer from certain limitations from today’s standard. For example, as a pioneering work, ROCO [24] was constructed long time ago with only 81k radiology images. MedICAT [29] contains 217k images, but are mostly consisted of compound figures. In this work, we tackle the above-mentioned limitations by introducing an automatic pipeline to generate dataset with subfigure-subcaption correspondence from scientific documents, consisting of three major stages: medical figure collection, subfigure separation, subcaption separation & alignment. The final dataset, PMC-OA, consisting of 1.65M image-text pairs, Fig. 1 and Fig. 3.

Along with the constructed dataset, we train a CLIP-style vision-language model for biomedical domain, termed as PMC-CLIP. The model is trained on PMC-OA with standard image-text contrastive (ITC) loss, and to encourage the joint interaction of image and text, masked language modeling (MLM) is also applied. We evaluate the pre-trained model on several downstream tasks, including medical image-text retrieval, medical image classification, and medical visual question answering (VQA). PMC-CLIP achieves state-of-the-art performance on various downstream tasks, surpassing previous methods significantly.

Overall, in this paper, we make the following contributions: **First**, we propose an automatic pipeline to construct high-quality image-text biomedical datasets from scientific papers, and construct an image-caption dataset via the proposed pipeline, named PMC-OA, which is $8\times$ larger than before. With the proposed pipeline, the dataset can be continuously updated. **Second**, we pre-train a vision-language model, termed as PMC-CLIP, on the constructed image-caption dataset, to serve as a foundation model for the biomedical domain. **Third**, we conduct thorough experiments on various tasks (retrieval, classification, and VQA), and obtain SOTA performance on most downstream datasets, demonstrating the superiority of PMC-OA and the potential of the foundation model PMC-CLIP. The dataset and pre-trained model will be made available to the community.

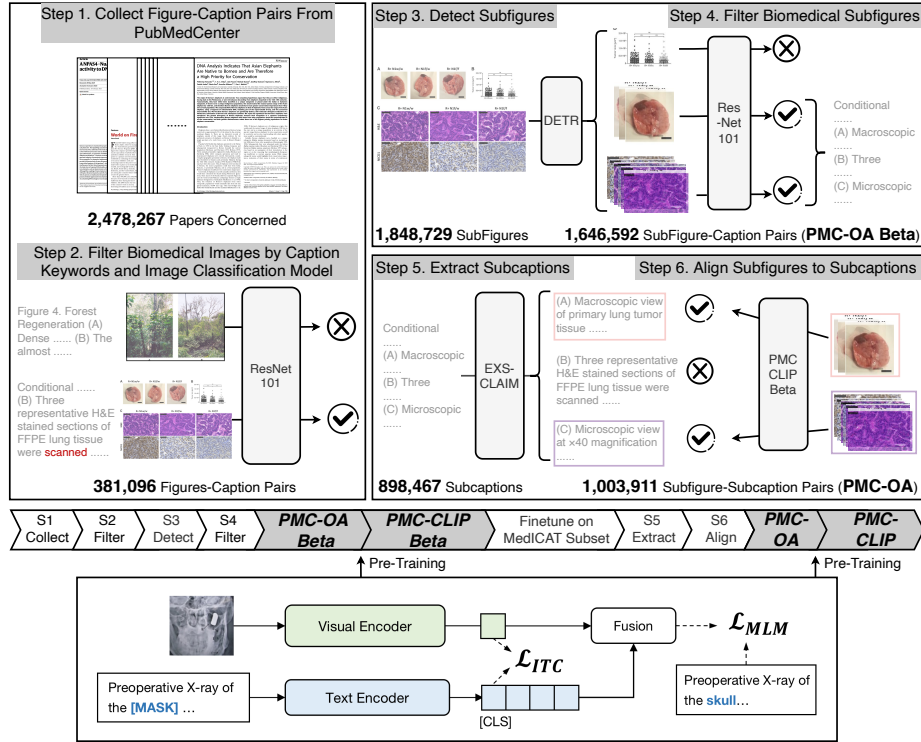


Fig. 2. The proposed pipeline to collect PMC-OA (upper) and the architecture of PMC-CLIP (bottom).

2 The PMC-OA Dataset

In this section, we start by describing the dataset collection procedure in Sec. 2.1, followed by a brief overview of PMC-OA in Sec. 2.2.

2.1 Dataset Collection

In this section, we detail the propose pipeline for creating PMC-OA, a large-scale dataset that contains 1.65M image-text pairs. The whole procedure consists of three major stages: (i) medical figure collection, (ii) subfigure separation, (iii) subcaption separation & alignment, as summarised in Fig. 2.

Medical Figure Collection (Step 1&2 in Fig. 2). We first extract figures and captions from PubMed Central, by the time of 2022-09-16, 2,478,267 available papers are covered and 12,211,907 figure-caption pairs are extracted. To derive medical figures, we follow the same procedure as in [29]: *first*, a set of medical keywords are pre-defined to filter the captions, deleting the figure-caption pair

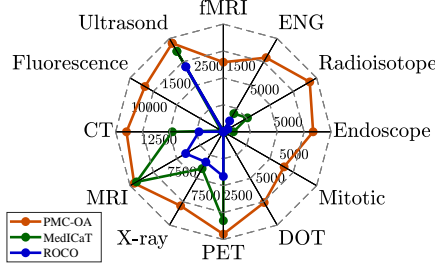
Follow Class T060 “*Diagnostic Procedure*” defined in UMLS [1]

with no keywords appearing in the caption; *second*, a ResNet-101 [12] trained on DocFigure [13] for scientific figure classification is applied to classify the remaining figures into 28 categories. We keep the figures with “Medical” class prediction scores in Top 4 of 28 categories, ending up with 381,096 medical figures.

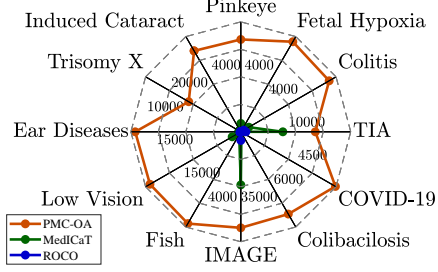
Subfigure Separation (Step 3&4 in Fig. 2). We randomly check around 300 figures from previous step, and find that around 80% of them are compound figures, i.e. multiple pannels. Here, we train a detector to break the compound figures into subfigures, specifically, we use a ResNet-34-based DETR model [3] with 4 encoder layers, 4 decoder layers, and 32 learnable queries. The detector is trained on the MedICaT subset [29], that has manually separated 2,069 compound figures and corresponding captions, which will be referred as MedICaTSub. For hyper-parameter tuning, we split the dataset into train and test set with a ratio of 3:1, our model obtains mAP@0.5 of 0.94 on the test set. To balance precision (0.93) and recall (0.94), we take the confidence threshold as 0.7. After breaking the compound figures, non-medical subfigures like charts may be mixed with medical ones, we therefore filter the derived subfigures with the classification model (repeat the first stage). Till here, we have obtained 1,646,592 subfigures from 378,717 compound figures, the number of captions is same as compound figures, thus each caption is assigned to an average of 4.3 subfigures at the moment. We termed this dataset as a **PMC-OA-Beta** version.

Subcaption Separation & Alignment (Step 5&6 in Fig. 2). To further align subfigure to its corresponding part within the full caption, i.e., subcaption, we need to break the captions into subcaptions first. Here, we apply the caption distributor provided in [28] to all the captions in PMC-OA Beta. **Note that**, the tool sometimes fail to separate the caption, as a result, we get 898,467 subcaptions from 239,905 separable captions, corresponding to 1,003,911 subfigures. To align each subfigure to the most related subcaption, we pretrain a CLIP-style model on **PMC-OA-Beta** (the training detail will be described in Sect. 3), then finetune it on MedICaTSub. Specifically, for each subfigure, we finetune the pre-trained model to match its subcaption with contrastive learning. We split the subset into train and test by 3:1, and the finetuned model achieves alignment accuracy=73% on test set. We finally align 1,003,911 subfigure-subcaption pairs. Along with the remaining 642,681 subfigure-caption pairs, we termed this dataset as **PMC-OA**. We consequently pretrain the **PMC-CLIP** on it.

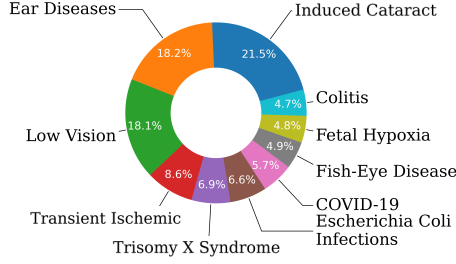
Discussion. Some pioneering works [24,29] have explored publicly available scientific documents to alleviate the data scarcity in the biomedical domain. As they provide valuable datasets for the community, they suffer from some limitations. Our work aims to improve the collection pipeline for a more massive, diverse, and accurate dataset: *First*, PMC-OA covers a wider range of papers (2,478,267) than ROCO [24](1,828,575) and MedICaT [29](131,410), and thus enlarge the dataset(1.6M). *Second*, different from ROCO [24], we maintain the non-radiology images, which makes PMC-OA contain more diverse biomedical data. We present the quantitative comparison in Fig. 3. *Third*, to the best of our knowledge, we are the first to integrate subfigures separation, subcaptions separation and the alignment into the data collection pipeline, which explicitly enlarges our dataset (ap-



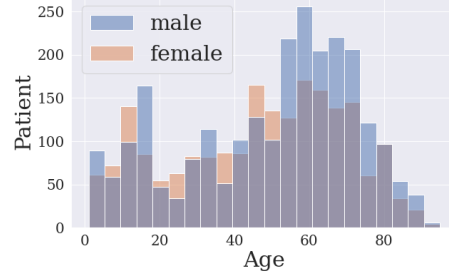
(a) Diagnostic procedure.



(b) Disease and findings.



(c) Disease distribution in PMC-OA.



(d) Patients' age & gender in PMC-OA.

Fig. 3. Statistical overview of PMC-OA.

proximately 8 times of MedICaT and 20 times of ROCO) while reducing the noise as much as possible.

2.2 Dataset Overview

In this section, we provide a brief statistical overview of the collected dataset PMC-OA from three different perspectives, *i.e.*, the diagnostic procedure, diseases and findings, and fairness.

Diagnostic Procedure. As shown in Fig. 3(a), PMC-OA covers a wide range of diagnostic procedures, spanning from common (*CT*, *MRI*, *X-ray*) to rare ones (*mitotic figure*), which is much diverse than before.

Disease and Findings. In PMC-OA, diseases are given in the free-form text, allowing for elaborate identification and analysis. For instance, eye diseases can be further categorized into Cataracts, Conjunctivitis, Macular degeneration, fish eye disease, etc. Fig. 3(b) illustrates the frequently used words in diagnosis and Fig. 3(c) shows their distribution.

Fairness. We provide the sex-ratio across ages in Fig. 3(d), as we can see PMC-OA is approximately gender-balanced, with 54% males.

Discussion. The detailed dataset statistics indicate the superiority of PMC-OA from three aspects: (i) diagnostic-procedure diversity; (ii) disease covering; (iii) population fairness. *First*, until now, the most widely-used text-image dataset is MIMIC-CXR[14] which contains only chest X-ray images, greatly limiting the potential of VLP methods and our dataset can compensate this. *Second*, diagnosis is always a crucial procedure in clinical, and the wide disease coverage in our dataset supports learning the shared patterns of diseases, promoting accurate auto-diagnosis. *Third*, the fairness on population ensures our dataset slightly suffers from patient characteristic bias, thus providing greater cross-center generalize ability.

3 Visual-language Pre-training

With our constructed image-caption dataset, we further train a visual-language model, termed as PMC-CLIP as shown in Fig. 2 (bottom). We describe the architecture in Sec. 3.1 and introduce the training objectives in Sec. 3.2

3.1 Architecture

Given a dataset with N image-caption pairs, *i.e.*, $\mathcal{D} = \{(\mathcal{I}_1, \mathcal{T}_1), \dots, (\mathcal{I}_N, \mathcal{T}_N)\}$, where $\mathcal{I}_i \in \mathbb{R}^{H \times W \times C}$ represents images, H, W, C are height, width, channel, and \mathcal{T}_i represents the paired text. We aim to train a CLIP-style visual-language model with an image encoder Φ_{visual} and a text encoder Φ_{text} .

In detail, given a specific image-caption pair $(\mathcal{I}, \mathcal{T})$ we encode it separately with a ResNet-based Φ_{visual} and a BERT-based Φ_{text} , the embedding dimension is denoted as d and the text token length as l :

$$\mathbf{v} = \Phi_{\text{visual}}(\mathcal{I}) \in \mathbb{R}^d, \quad (1)$$

$$\mathbf{T} = \Phi_{\text{text}}(\mathcal{T}) \in \mathbb{R}^{l \times d}, \quad \mathbf{t} = \mathbf{T}_0 \in \mathbb{R}^d, \quad (2)$$

where \mathbf{v} represents the embedding for the whole image, \mathbf{T} refers to the sentence embedding, and \mathbf{t} denotes the embedding for [CLS] token.

3.2 Training Objectives

In this section, we train the visual-language model with two objectives, *i.e.*, Image-Text Contrastive learning and Masked Language Modeling.

Image-Text Contrastive Learning (ITC). We implement ITC loss following CLIP [25], that aims to match the corresponding visual and text representations from one sample. In detail, denote batch size as b , we calculate the softmax-normalized cross-modality dot product similarity between the current visual/text embedding (\mathbf{v} / \mathbf{t}) and all samples within the batch, termed as $p^{\text{i}2\text{t}}, p^{\text{t}2\text{i}} \in \mathbb{R}^b$, and the final ITC loss is:

$$L_{\text{ITC}} = \mathbb{E}_{(\mathcal{I}, \mathcal{T}) \sim \mathcal{D}} [\text{CE}(y^{\text{i}2\text{t}}, p^{\text{i}2\text{t}}) + \text{CE}(y^{\text{t}2\text{i}}, p^{\text{t}2\text{i}})], \quad (3)$$

where y^{i2t}, y^{t2i} refer to one-hot matching labels, CE refers to InfoNCE loss [22].

Masked Language Modeling (MLM). We implement MLM loss following BERT [7]. The network is trained to reconstruct the masked tokens from context contents and visual cues. We randomly mask the word in texts with a probability of 15% and replace it with Token [MASK]. We concatenate the image embedding \mathbf{v} with the sentence embedding \mathbf{T} , input it into a self-attention transformer-based fusion module Φ_{fusion} , and get the prediction for the masked token at the corresponding position in the output sequence, termed as $p^{\text{mask}} = \Phi_{\text{fusion}}(\mathbf{v}, \mathbf{T})$. Let y^{mask} denote the ground truth, and the MLM loss is:

$$L_{\text{MLM}} = \mathbb{E}_{(\mathcal{I}, \mathcal{T}) \sim \mathcal{D}} [\text{CE}(y^{\text{mask}}, p^{\text{mask}})] \quad (4)$$

Total Training Loss. The final loss is the weighted sum of the above two:

$$L = L_{\text{ITC}} + \lambda L_{\text{MLM}}, \quad (5)$$

where λ is a hyper-parameter deciding the weight of L_{MLM} , set as 0.5 by default.

4 Experiment Settings

Here, we start by introducing the compared datasets in Sec. 4.1, describe the downstream tasks in Sec. 4.2, cover the implementation details in Sec. 4.3.

4.1 Pre-training Datasets

ROCO [24] is a image-caption dataset collected from PubMed [27]. It filters out all the compound or non-radiological images, and consists of 81K samples.

MedICaT [29] extends ROCO to 217K samples (image-caption pairs), however, 75% of its figures are compound ones, *i.e.* one figure with multiple subfigures.

MIMIC-CXR [14] is the largest chest X-ray dataset, containing 377,110 samples (image-report pairs). Each image is paired with a clinical report describing findings from doctors.

4.2 Downstream Tasks

Image-Text Retrieval (ITR). ITR contains both image-to-text(I2T) and text-to-image(T2I) retrieval. We train PMC-CLIP on different datasets, and evaluate on the ROCO testset. **Note that**, we have explicitly conducted duplication between our data and ROCO, our reported results thus resembles *zero-shot* evaluation. Following previous works [29, 4, 5], we sample 2,000 image-text pairs from ROCO’s testset for evaluation.

Classification. We finetune the model for different downstream tasks that focus on image classification. Specifically, MedMINIST [31] contains 12 tasks in total, and it covers primary data modalities in biomedical images, including Colon Pathology, Dermatoscope, Retinal OCT, etc.

Visual Question Answering (VQA). We evaluate on the official dataset split of VQA-RAD [16], SLAKE [19], where SLAKE is composed of 642 images and 14,028 questions and VQA-RAD contains 315 images and 3,515 questions. The questions in VQA-RAD and Slake are categorized as close-ended if answer choices are limited, otherwise open-ended. While adapting our PMC-CLIP to VQA task, we maintain most of the pre-trained parameters, including visual encoder, text encoder, and fusion. Specifically, the image and question are fed into PMC-CLIP, the output embedding from fusion module is then used to compute similarity between each of the answer candidates (also encoded with text encoder). To strive for better adaptation, we use 10 learnable answer and question prompt vectors in text encoder respectively.

Metrics. We use Accuracy and AUC are for classification, Recall@K(K=1,5,10) for retrieval and Accuracy for VQA.

4.3 Implementation Details

For the visual and text encoders, we adopt ResNet50 [12] and PubmedBERT [11]. And we use 4 transformer layers for the fusion module. For input data, we resize each image to 224×224 . During pre-training, our text encoder is initialized from PubmedBERT, while the vision encoder and fusion module are trained from scratch. We use AdamW [20] optimizer with $lr = 1 \times 10^{-4}$. We train on GeForce RTX 3090 GPUs with batch size 128 for 100 epochs. The first 10 epochs are set for warming up.

5 Result

We conduct experiments to validate our proposed dataset, and the effectiveness of models trained on it. In Sec. 5.1, we first compare with existing large-scale biomedical datasets on the image-text retrieval task to demonstrate the superiority of PMC-OA. In Sec. 5.2, we finetune the model (pre-trained on PMC-OA) across three different downstream tasks, namely, retrieval, classification, and visual question answering. And we also perform a thorough empirical study of the pretraining objectives and the model architectures in Sec. 5.3. **Note that**, for all experiments, unless specified otherwise, we use the default setting: ResNet50 for image encoder, and pre-train with both ITC and MLM objectives.

5.1 PMC-OA surpasses SOTA large-scale biomedical dataset

As shown in Tab 1, we pre-train PMC-CLIP on different datasets and evaluate the retrieval on ROCO test set. The performance is largely improved while simply switching to our dataset, confirming the significance of it.

Table 1. Ablation studies on pre-training dataset.

Methods	Pretrain Data	DataSize	I2T			T2I		
			R@1	R@5	R@10	R@1	R@5	R@10
PMC-CLIP	ROCO	173 K	12.30	35.28	46.52	13.36	35.84	47.38
PMC-CLIP	MedICaT	173 K	12.30	35.28	46.52	13.36	35.84	47.38
PMC-CLIP	PMC-OA Beta	1.6 M	30.42	59.11	70.16	27.92	55.99	66.35
PMC-CLIP	PMC-OA	1.6 M	31.41	61.15	71.88	28.02	58.33	69.69

5.2 PMC-CLIP achieves SOTA across downstream tasks

To evaluate the learnt representation in PMC-CLIP, we compare it with several state-of-the-art approaches across various downstream tasks, including image-text retrieval, image classification, and visual question answering.

Image-Text Retrieval. As shown in Tab. 2, we report a state-of-the-art result on image-text retrieval. On I2T Rank@10, PMC-CLIP outperforms previous state-of-the-art by 8.1%. It is worth mentioning that, the training set of ROCO has been used during pretraining in M3AE [4], ARL [5]. While our dataset does not contain data from ROCO.

Table 2. Zero-shot Image-Text Retrieval on ROCO. Dark and light grey colors highlight the top and second best results on each metric.

Methods	Pretrain Data	DataSize	I2T			T2I		
			R@1	R@5	R@10	R@1	R@5	R@10
ViLT [15]	COCO, VG, SBU, GCC	4.1M	11.90	31.90	43.20	9.75	28.95	41.40
METER [9]	COCO, VG, SBU, GCC	4.1M	14.45	33.30	45.10	11.30	27.25	39.60
M3AE [4]	ROCO, MedICaT	233 K	19.10	45.60	61.20	19.05	47.75	61.35
ARL [5]	ROCO, MedICaT, CXR	233 K	23.45	50.60	62.05	23.50	49.05	63.00
PMC-CLIP	PMC-OA	1.6 M	31.41	61.15	71.88	28.02	58.33	69.69

Image Classification. To demonstrate the excellent transferability of PMC-CLIP, we finetune it on MedMNIST and compare it with SOTA methods, *i.e.*, DWT-CV [6] and SADA [10]. We present the results of 3 of 12 sub-tests here, and the full results can be found in the supplementary material. As shown in Tab. 3, PMC-CLIP obtains consistently higher results, and it is notable that finetuning from PMC-CLIP achieves significant performance gains compared with training from scratch with ResNet.

Table 3. Classification results on MedMNIST. Dark and light grey colors highlight the top and second best results on each metric.

Dataset Method	PneumoniaMNIST		BreastMNIST		DermaMNIST	
	AUC↑	ACC↑	AUC↑	ACC↑	AUC↑	ACC↑
ResNet50 [12]	96.20	88.40	91.90	86.60	91.20	73.10
DWT-CV [6]	95.69	88.67	89.77	85.68	91.67	74.75
SADA [10]	98.30	91.80	91.50	87.80	92.70	75.90
PMC-CLIP	99.02	95.35	94.56	91.35	93.41	79.80

Visual Question Answering. VQA requires model to learn finer grain visual and language representations. As Table 4 shows, we surpass SOTA method M3AE in 5 out of 6 results.

Table 4. VQA results on VQA-RAD and Slake

Methods	VQA-RAD			Slake		
	Open	Closed	Overall	Open	Closed	Overall
MEVF-BAN [21]	49.20	77.20	66.10	77.80	79.80	78.60
CPRD-BAN [18]	52.50	77.90	67.80	79.50	83.40	81.10
M3AE [4]	67.23	83.46	77.01	80.31	87.82	83.25
PMC-CLIP	67.00	84.00	77.60	81.90	88.00	84.30

5.3 Ablation Study

For the effectiveness illustration of the pretraining objectives (*ITC*, *MLM*), we evaluate PMC-CLIP pre-trained on PMC-OA for ablation studies. To validate the necessity of fusion module, we compare MLM using text-only context (*MLM-T*), and MLM with visual cues (*MLM-V*).

As shown in Tab. 5, we can draw the following observations: *First*, ITC objective is essential for pretraining, and contributes most of the performance (ID 1). *Second*, MLM using only text context works as a regularization term (ID 2). *Third*, With incorporation of visual features, the model learns finer grain correlation between image-caption pairs, and achieve the best results (ID 3).

Table 5. Ablation studies of pre-training objectives

ID	Methods			ROCO			PneumoniaMNIST		BreastMNIST		DermaMNIST	
	ITC	MLM-T	MLM-V	R@1	R@5	R@10	AUC	ACC	AUC	ACC	AUC	ACC
1	✓			25.53	53.48	64.64	98.07	93.56	93.04	89.74	92.77	77.91
2	✓	✓		29.17	57.55	68.26	98.53	94.31	94.39	90.71	92.76	78.00
3	✓		✓	29.72	59.74	70.79	99.02	95.35	94.56	91.35	93.41	79.80

6 Conclusion

In this paper, we present a large-scale dataset in biomedical domain, named PMC-OA, by collecting image-caption pairs from abundant scientific documents. We train a CLIP-style model on on PMC-OA, termed as PMC-CLIP, it achieves SOTA performance across various downstream biomedical tasks, including image-text retrieval, image classification, visual question answering. With the automatic collection pipeline, we believe the dataset will benefit the research community, fostering the development of foundation models in biomedical domain.

References

1. Olivier Bodenreider. The unified medical language system (umls): integrating biomedical terminology. *Nucleic acids research*, 32(suppl_1):D267–D270, 2004. [3](#)
2. Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. [1](#)
3. Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 213–229. Springer, 2020. [4](#)
4. 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. [7](#), [9](#), [10](#)
5. 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. [7](#), [9](#)
6. Jianhong Cheng, Hulin Kuang, Qichang Zhao, Yahui Wang, Lei Xu, Jin Liu, and Jianxin Wang. Dwt-cv: Dense weight transfer-based cross validation strategy for model selection in biomedical data analysis. *Future Generation Computer Systems*, 135:20–29, 2022. [9](#)
7. Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. [7](#)
8. Ming Ding, Wendi Zheng, Wenyi Hong, and Jie Tang. Cogview2: Faster and better text-to-image generation via hierarchical transformers. *arXiv preprint arXiv:2204.14217*, 2022. [1](#)
9. Zi-Yi Dou, Yichong Xu, Zhe Gan, Jianfeng Wang, Shuohang Wang, Lijuan Wang, Chenguang Zhu, Pengchuan Zhang, Lu Yuan, Nanyun Peng, et al. An empirical study of training end-to-end vision-and-language transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18166–18176, 2022. [9](#)
10. Xiaolong Ge, Yanpeng Qu, Changjing Shang, Longzhi Yang, and Qiang Shen. A self-adaptive discriminative autoencoder for medical applications. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(12):8875–8886, 2022. [9](#)
11. Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23, 2021. [8](#)
12. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [4](#), [8](#), [9](#)
13. KV Jobin, Ajoy Mondal, and CV Jawahar. Docfigure: A dataset for scientific document figure classification. In *2019 International Conference on Document Analysis and Recognition Workshops (ICDARW)*, volume 1, pages 74–79. IEEE, 2019. [4](#)

14. Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Roger G Mark, and Steven Horng. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data*, 6(1):317, 2019. [6](#), [7](#)
15. 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. [9](#)
16. 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. [8](#)
17. 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*, 34:9694–9705, 2021. [1](#)
18. 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 *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part II 24*, pages 210–220. Springer, 2021. [10](#)
19. 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. [8](#)
20. Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. [8](#)
21. 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 *Medical Image Computing and Computer Assisted Intervention–MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part IV 22*, pages 522–530. Springer, 2019. [10](#)
22. Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. [7](#)
23. Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022. [1](#)
24. 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: 7th Joint International Workshop, CVII-STENT 2018 and Third International Workshop, LABELS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16, 2018, Proceedings 3*, pages 180–189. Springer, 2018. [2](#), [4](#), [7](#)
25. 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. [1](#), [6](#)
26. Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. [1](#)
27. Richard J Roberts. Pubmed central: The genbank of the published literature, 2001. [2](#), [7](#)

28. Eric Schwenker, Weixin Jiang, Trevor Spreadbury, Nicola Ferrier, Oliver Cossairt, and Maria KY Chan. Exsclaim!—an automated pipeline for the construction of labeled materials imaging datasets from literature. *arXiv preprint arXiv:2103.10631*, 2021. [4](#)
29. Sanjay Subramanian, Lucy Lu Wang, Sachin Mehta, Ben Bogin, Madeleine van Zuylen, Sravanthi Parasa, Sameer Singh, Matt Gardner, and Hannaneh Hajishirzi. Mediat: A dataset of medical images, captions, and textual references. *arXiv preprint arXiv:2010.06000*, 2020. [2](#), [3](#), [4](#), [7](#)
30. 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. [1](#)
31. Jiancheng Yang, Rui Shi, Donglai Wei, Zequan Liu, Lin Zhao, Bilian Ke, Hanspeter Pfister, and Bingbing Ni. Medmnist v2—a large-scale lightweight benchmark for 2d and 3d biomedical image classification. *Scientific Data*, 10(1):41, 2023. [7](#)