# Does CLIP Benefit Visual Question Answering in the Medical Domain as Much as it Does in the General Domain?

**Sedigheh Eslami**

D4L data4life gGmbH

Potsdam, Germany

**Gerard de Melo**

Hasso Plattner Institute

Potsdam, Germany

**Christoph Meinel**

Hasso Plattner Institute

Potsdam, Germany

sedigheh.eslami@data4life.care

{gerard.demelo, christoph.meinel}@hpi.de

December 30, 2021

## ABSTRACT

**Objective:** Contrastive Language–Image Pre-training (CLIP) has shown remarkable success in learning with cross-modal supervision from extensive amounts of image–text pairs collected online. Thus far, the effectiveness of CLIP has been investigated primarily in general-domain multimodal problems. This work evaluates the effectiveness of CLIP for the task of Medical Visual Question Answering (MedVQA). To this end, we present PubMedCLIP, a fine-tuned version of CLIP for the medical domain based on PubMed articles.
**Materials and Methods:** Our experiments are conducted on two MedVQA benchmark datasets and investigate two MedVQA methods, MEVF (Mixture of Enhanced Visual Features) and QCR (Question answering via Conditional Reasoning). For each of these, we assess the merits of visual representation learning using PubMedCLIP, the original CLIP, and state-of-the-art MAML (Model-Agnostic Meta-Learning) networks pre-trained only on visual data. We open source the code for our MedVQA pipeline and pre-training PubMedCLIP.
**Results:** CLIP and PubMedCLIP achieve improvements in comparison to MAML's visual encoder. PubMedCLIP achieves the best results with gains in the overall accuracy of up to 3%. Individual examples illustrate the strengths of PubMedCLIP in comparison to the previously widely used MAML networks.
**Discussion and conclusion:** Visual representation learning with language supervision in PubMed-CLIP leads to noticeable improvements for MedVQA. Our experiments reveal distributional differences in the two MedVQA benchmark datasets that have not been imparted in previous work and cause different back-end visual encoders in PubMedCLIP to exhibit different behavior on these datasets. Moreover, we witness fundamental performance differences of VQA in general versus medical domains.

*Keywords* Medical visual question answering · Deep representation learning · CLIP · PubMedCLIP

## 1 BACKGROUND AND SIGNIFICANCE

Medical visual question answering (MedVQA) is the task of answering natural language questions about a given medical image. To solve such multimodal tasks, a system must interpret both visual and textual data as well as infer the associations between a given image and a pertinent question sufficiently well to elicit an answer Antol et al. [2015]. The development of MedVQA has considerable potential to benefit healthcare systems, as it may aid clinicians in interpreting medical images, obtaining more accurate diagnoses by consulting a second opinion, and ultimately, may expedite and improve patient care. Achieving this in the medical domain in particular is non-trivial, as we suffer from a

general lack of sufficient and balanced training data. The ImageCLEF community hosts annual MedVQA challenges Abacha et al. [2020], Ben Abacha et al. [2021], where new VQA datasets using PubMed articles are released. However, there are concerns about whether the question–answer pairs in these datasets are realistic and clinically relevant Lau et al. [2018]. For example, in the VQA-Med 2021 Challenge Ben Abacha et al. [2021], the dataset consisted entirely of questions asking about the category of abnormality in the image. This lack of diversity in the semantics of the questions meant that the winning teams were able to treat the MedVQA problem as a multi-class image classification task, without any need to interpret the questions Gong et al. [2021], Eslami et al. [2021]. Lau *et al.*Lau et al. [2018] published VQA-RAD as the first public benchmark dataset comprising realistic and clinically relevant question–answer pairs generated by expert clinicians and radiologists. Recently, Liu *et al.*Liu et al. [2021a] created the bilingual SLAKE dataset that includes not only clinically relevant data, but also mask and bounding box annotations for images, which are beneficial for semantic segmentation and detection of organs in medical images.

Current approaches for this multimodal task adopt deep neural encoders to interpret the image and the question and then pick a corresponding answer. They typically consist of four main components: a visual encoder, question encoder, attention-based fusion of vision and text features, and an answer classifier Vu et al. [2020], Zhan et al. [2020], Nguyen et al. [2019], Pan et al. [2021], Liu et al. [2021b]. Skip-thought vectors, LSTM, and GRU recurrent neural networks have been popular question encoders in prior work. Bilinear attention networks Kim et al. [2018], stacked attention networks Yang et al. [2016], and element-wise production are popular as multimodal pooling approaches in MedVQA. With regard to the visual encoder, the majority of previous MedVQA papers Pan et al. [2021], Zhan et al. [2020], Nguyen et al. [2019] employ the Mixture of Enhanced Visual Features (MEVF) Nguyen et al. [2019]. It consists of two modules: 1. the pre-trained meta learning module, which uses Model-Agnostic Meta-Learning (MAML) Finn et al. [2017] with the objective of solving a $k$-shot $n$-way classification problem with the abnormality status of different organs as classes, 2. the Convolutional Denoising Autoencoder (CDAE) Masci et al. [2011] module in order to denoise the medical image. However, the pre-training of MEVF is custom-tailored for the particular challenges encountered in the VQA-RAD benchmark dataset and is specifically designed for the organs present in this dataset, i.e., the chest, brain, and abdomen, limiting its generalizability to other settings. Liu *et al.*Liu et al. [2021b] similarly restricted the objective of their visual encoding to chest, brain, and abdomen, and pre-train three separate visual feature extraction teacher models for these respective body regions. Furthermore, they distilled the three teacher models into a smaller student model by contrastive representation distillation. This motivated us to design an alternative model, PubMedCLIP, which learns features in medical images of various body organs and is not limited to only a few regions.

Transfer learning and making use of pre-trained models has become an inseparable part of representation learning in computer vision and natural language processing. Recent work Radford et al. [2021], Cho et al. [2021], Su et al. [2019] has shown improvements of visual and textual encoders when learning from the contrast of image–text pairs and using natural language as supervision in addition to just visual images. This trend of improvements has also been observed in various classification use cases in the medical domain Zhang et al. [2020]. Among these approaches, the contrastive pre-training of language–image data in CLIP Radford et al. [2021] has been particularly successful. CLIP was trained using a very large number of image–text pairs acquired from the Internet with close to zero additional human annotation. This aspect is particularly useful for the medical domain, since data annotation requires expert medical knowledge and thus is often an expensive and time-consuming obstacle. Following CLIP, we investigate to what extent using publicly available medical image–text pairs without any further annotation can be useful for the MedVQA task. To this end, we consider a large number of medical image–text pairs obtained from PubMed articles and use them to train PubMedCLIP. We further examine the outcomes when incorporating PubMedCLIP as the visual encoder into state-of-the-art MedVQA methods. We investigate whether fine-tuning CLIP for the MedVQA task benefits medical VQA as much as it benefits general-domain VQA, as observed in previous work Shen et al. [2021].

To the best of our knowledge, this is the first study introducing a PubMed-optimized CLIP and assessing the effects of using CLIP for MedVQA. In contrast to previous visual encoders used in MedVQA, PubMedCLIP is pre-trained using medical images from a diverse range of body regions and is not restricted to only brain, chest, and abdomen images. We conduct extensive experiments on two MedVQA benchmark datasets and employ diverse back-end visual encoders in PubMedCLIP. Our experiments reveal that using PubMedCLIP as a pre-trained visual encoder improves previous models by up to $3\%$.

## 2 MATERIALS AND METHODS

### 2.1 PubMedCLIP

Our first step was to consider the original CLIP, which has been pre-trained on general-domain images encountered online, and fine-tune it using medical image–text pairs. To this end, we drew on the Radiology Objects in COntext (ROCO) dataset Pelka et al. [2018], which provides over $80K$ samples. ROCO includes diverse imaging modalities
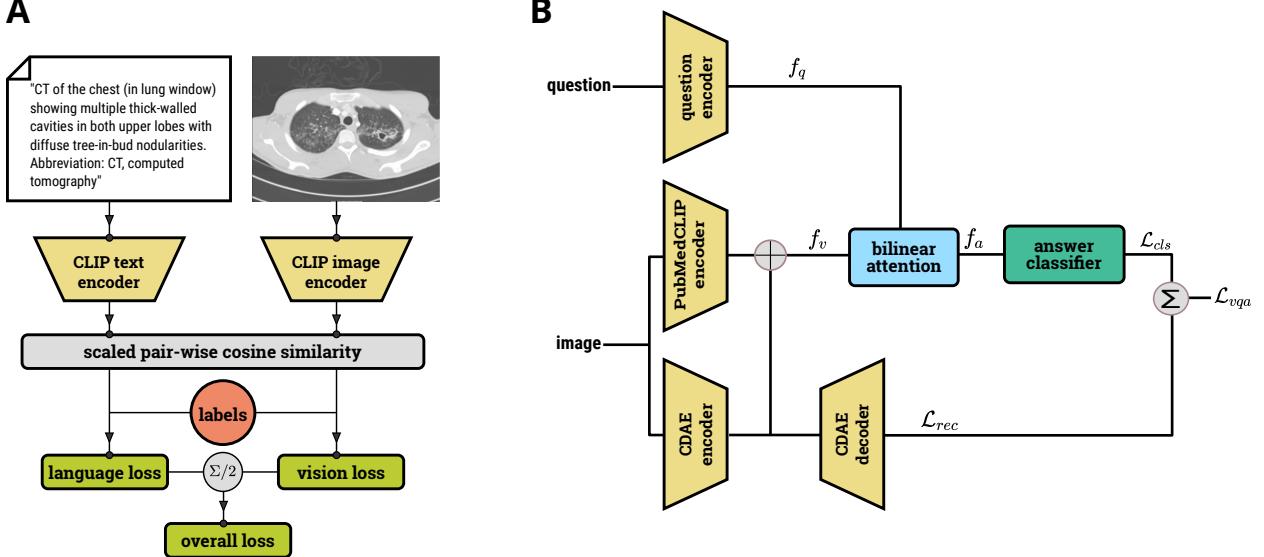
**A**



**B**



**Figure 1:** (A) Overview of how PubMedCLIP is pre-trained. (B) Schematic of MedVQA backbone with PubMedCLIP pre-trained visual encoder.

such as ultrasound, X-Ray, fluoroscopy, PET scans, mammography, MRI, angiography, from various human body regions, e.g., head, neck, jaw and teeth, spine, chest, abdomen, hand, foot, knee, and pelvis. The image–text pairs in this dataset are captured from PubMed articles. The texts here are taken from the relatively short captions (average length of 20 words) associated with images in the articles, which provide rich explanatory information about the content of images. To the best of our knowledge, ROCO is the only large-scale publicly available medical dataset that includes image–text pairs for a diverse range of body organs and imaging modalities.

In this work, the training and validation data splits from the original paper Pelka et al. [2018] were used to fine-tune CLIP for the medical domain, with ViT32 Vision Transformer Dosovitskiy et al. [2021], ResNet RN-50 and RN-50x4 He et al. [2016] visual encoder back-ends. With respect to the maximum text length accepted by CLIP, which is 76, we trimmed any longer captions, while zero-padding shorter ones.

We refer to the resulting fine-tuned model as PubMedCLIP. PubMedCLIP was trained for $50$ epochs with a batch size of $64$, and Adam optimization Kingma and Ba [2014] with a learning rate of $10^{-5}$. The source code along with further implementation details can be found at: https://github.com/sarahESL/PubMedCLIP. An overview of the training procedure for PubMedCLIP is given in Figure 1 (A). Text and image are encoded separately using CLIP. The cosine similarity between text and image features is computed. Finally, the vision and language cross-entropy loss values are computed and their average is considered as the overall loss value.

## 2.2 PubMedCLIP in MedVQA

Our goal is to investigate the effect of using PubMedCLIP as a pre-trained visual encoder in MedVQA models. VQA in this work is considered as a classification problem, where the objective is to find a mapping function $f$ that maps an image–question pair $(v_i, q_i)$ to the natural language answer $a_i$. For our investigation, we considered two prominent MedVQA methods, MEVF Zhan et al. [2020] and QCR Nguyen et al. [2019], that adopt MEVF as their visual encoders. To assess the contribution of PubMedCLIP, we modified the MEVF by substituting its pre-trained MAML module with PubMedCLIP. Hence, the representative visual feature in our work is the concatenation of the output of the PubMedCLIP network and the CDAE encoder network. Both models use GloVe word embeddings Pennington et al. [2014] followed by an LSTM in order to encode questions. Furthermore, the multimodal pooling mechanism for combining question and image features is BAN Kim et al. [2018] in both models. We retained the same question encodings, multimodal fusion, and objective functions proposed for MEVF Nguyen et al. [2019] and QCR Zhan et al. [2020], respectively, and only replaced the visual encoders.

The overall objective is to minimize the error of answer classification and image reconstruction, denoted as

$$\mathcal{L}_{\text{vqa}} = \mathcal{L}_{\text{cls}} + \mathcal{L}_{\text{rec}}. \qquad (1)$$

3

Following previous work Teney et al. [2018], a sigmoid layer preceding a binary cross-entropy loss computation is used for the classification. The loss function for the autoencoder reconstruction is mean squared error. A schematic architecture of the backbone of our work is shown in Figure 1 (B). The answer classifier is a two-layer feed-forward network with the ReLU activation function, as proposed for BAN Kim et al. [2018].

### 2.3 Datasets

We conducted our experiments using two well-known datasets:

1. **VQA-RAD** Lau et al. [2018] consists of 315 medical images and 3,515 question–answer pairs. We follow previous work by adopting the data split proposed for MEVF Nguyen et al. [2019]. We notice that all the images in the test dataset are also present in the training set. However, the set of question–answer pairs for these images in the test set are unseen in the training set.

2. The **SLAKE** Liu et al. [2021a] dataset consists of English and Chinese questions. In this work, we utilize the English subset of the dataset, comprising 642 images and more than 7,000 question–answer pairs. Using the original data split, we observe that in contrast to VQA-RAD, all the images in the test set of SLAKE are unseen in the training set.

### 2.4 Experimental setup

In order to ensure a fair comparison, our experiments generally followed the same setups used in the MEVF and QCR studies. For both methods, Adam optimization was invoked for training. MEVF was trained for 20 epochs, QCR for 200 epochs. When using PubMedCLIP as the pre-trained visual encoder, we set the learning rate to $1 \times 10^{-3}$ and $2 \times 10^{-3}$ and the batch size to 16 and 32 in QCR and MEVF, respectively. All implementations are based on the PyTorch framework Paszke et al. [2019]. We ran the original MEVF and QCR on our machine and report the results here to have a fair comparison. Due to the non-deterministic behaviour of the cuDNN library used in CUDA convolution operations Pham et al. [2020], we observed non-deterministic results in different runs. For a more robust comparison, we repeated all experiments 10 times and report the average accuracy scores.

## 3 RESULTS

The results of our experiments are reported in Table 1. We provide the overall accuracy along with the accuracy of answering open-end and closed-end questions. It is observed that the performance of both MEVF and QCR approaches are improved when adopting CLIP and PubMedCLIP as the pre-trained visual encoder. Furthermore, results of PubMedCLIP show up to $1\%$ improvement in comparison to the original CLIP. For the VQA-RAD dataset, PubMedCLIP with the ResNet-50 backend achieves the best results, improving the overall accuracy of MEVF up to $6\%$ and for QCR up to $3\%$ percent. Results on the SLAKE dataset indicate that PubMedCLIP with the back-end ViT32 Vision Transformer visual encoder attains the best accuracy. It enhances MEVF up to $3\%$ and QCR up to $2\%$. We witness the same trend of improvement among overall, open-end, and closed-end accuracy scores.

## 4 DISCUSSION

The fact that ResNet-50 for VQA-RAD and ViT for SLAKE dataset achieve the best results suggests that there are underlying differences in the question type distribution in these datasets. As shown in Figure 2, the majority of the questions in the VQA-RAD dataset ask about the presence of an abnormality in the images. This requires the visual encoder to detect local features and local abnormalities in the image. Therefore, the CNN-based ResNet model with better visual localization outperforms the Vision Transformer. However, on the SLAKE dataset, the majority of the questions are from the type "organ", asking which organ is present in the image. For such cases, the visual encoder needs to be able to acquire a holistic overall understanding of the content of the image and thus also capture long-range dependencies of image patches. Vision Transformers indeed are capable of accounting for such features Yu et al. [2021], and hence perform better on the SLAKE dataset. Figure 2 plots the distribution of question types for the top 5 frequent types in both datasets.

In Figures 3 and 4, we provide examples from the VQA-RAD and SLAKE datasets, respectively. Our goal is to illustrate the performance of the original MEVF and QCR in comparison with QCR when PubMedCLIP is used as the visual encoder for the MedVQA task. We refer to the QCR model with PubMedCLIP simply as PubMedCLIP in these figures. Examples from both datasets in Figures 3 and 4 demonstrate that the MEVF model has difficulties correctly comprehending which organ is depicted in the image. For example, regardless of the asked question, in the
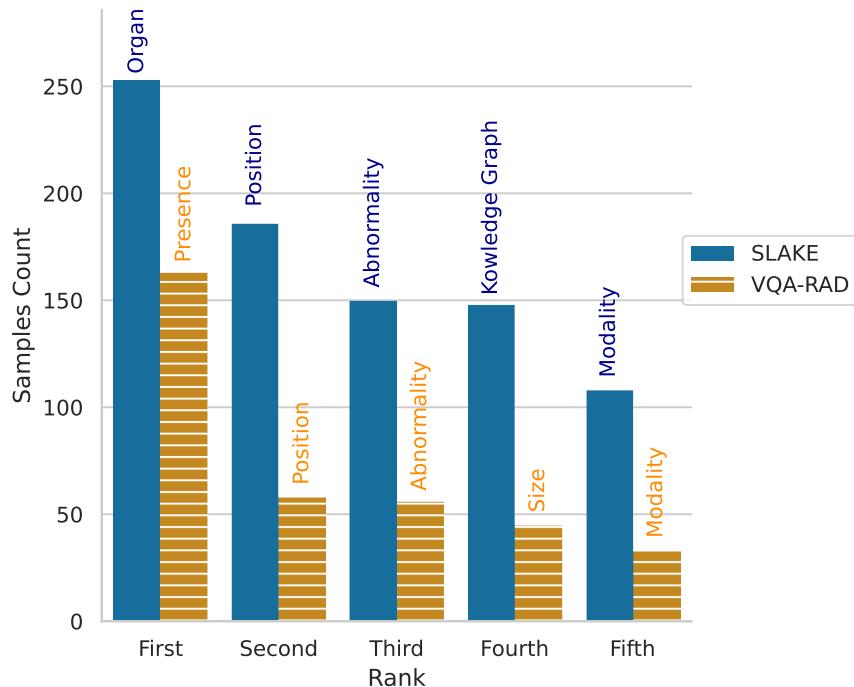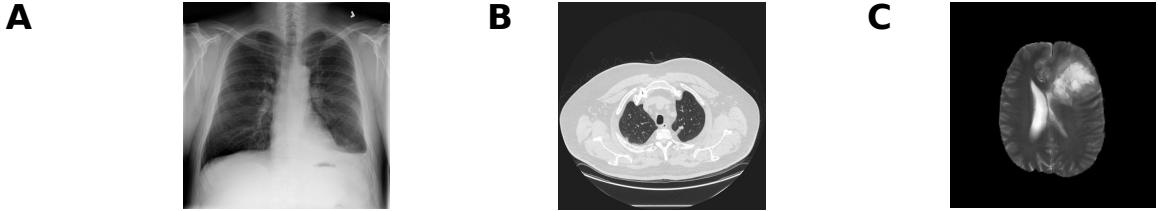
**Figure 2:** Distribution of the top 5 frequent question types in VQA-RAD and SLAKE datasets.



**Figure 3:** Examples from VQA-RAD dataset.

| MedVQA model | Visual encoder | VQA-RAD Accuracy | | | SLAKE Accuracy | | |
|---|---|---|---|---|---|---|---|
| | | open | closed | overall | open | closed | overall |
| MEVF | MAML + AE | 42.1% | 73.2% | 60.8% | 74.1% | 77.5% | 75.5% |
| | CLIP-ViT-B + AE | 50.8% | 75% | 65.4% | 75.8% | 80.5% | 77.7% |
| | CLIP-RN50 + AE | 47% | 77.4% | 65.4% | 75.7% | 79.6% | 77.2% |
| | CLIP-RN50x4 + AE | 46.8% | 76.6% | 64.8% | 75.9% | 79.1% | 77.2% |
| | PubMedCLIP-ViT-B + AE | 48.9% | 76.7% | 65.5% | **76.5%** | **80.4%** | **78%** |
| | PubMedCLIP-RN50 + AE | **48.6%** | **78.1%** | **66.5%** | 76.2% | 79.9% | 77.6% |
| | PubMedCLIP-RN50x4 + AE | 47.1% | 77.8% | 65.6% | 76.6% | 79.1% | 77.6% |
| QCR | MAML + AE | 56% | 77.9% | 69.2% | 76.8% | 80.6% | 78.3% |
| | CLIP-ViT-B + AE | 57.6% | 79.5% | 70.7% | 78.6% | 81% | 79.5% |
| | CLIP-RN50 + AE | 58.3% | 80% | 71.3% | 78.2% | 81.5% | 79.7% |
| | CLIP-RN50x4 + AE | 59.9% | 79.4% | 71.3% | 77.6% | 80.5 | 78.7 |
| | PubMedCLIP-ViT-B + AE | 58.4% | 79.5% | 71.1% | **78.4%** | **82.5%** | **80.1%** |
| | PubMedCLIP-RN50 + AE | **60.1%** | **80%** | **72.1%** | 77.8% | 81.4% | 79.3% |
| | PubMedCLIP-RN50x4 + AE | 60% | 79.7% | 71.8% | 77.7% | 81.3% | 79.1% |

**Table 1:** Accuracy scores on VQA-RAD and SLAKE datasets.

**A**

Question: Where does the image represent in the body?

Answer: chest
MEVF: atelectasis, effusion ✗
QCR: lower left lung ✗
PubMedCLIP: chest ✓

**B**

What diseases are included in the picture?

lung cancer
left lung ✗
none ✗
lung cancer ✓

**C**

Where is the brain non-enhancing tumor?

upper left lobe
rectum, small bowel ✗
right lung ✗
upper left lobe ✓

**Figure 4:** Examples from SLAKE dataset.

left image in Figure 3, we observe that although the given image is a chest x-ray, the answer that MEVF provides is related to the abdominal region. This behaviour is seen for the right examples as well, where the given image is from the brain, but the predicted answer from MEVF relates to the chest. The same behaviour is further also observed in Figure 4. From this perspective, QCR appears to be providing answers that are at least relevant to the given image and question, although it fails to select the correct answer. As an example, for the left image in Figure 3, QCR understands the question as well as the image, but cannot provide the coarse correct answer. In the left image from Figure 4, QCR appears to understand that it relates to the lung/chest area, but it has apparent difficulties in comprehending the question and providing the correct answer. In contrast, using QCR with PubMedCLIP shows an improvement and results in providing answers that are correct throughout all examples in Figures 3 and 4.

**A**

**B**

**C**

| Question: | Are there multiple or just 1 metastatic focus? | What are the hyperdensities on the periphery of the image? | What is the biological sex of the patient? |
|---|---|---|---|
| Answer: | one | ribs | female |
| MEVF: | right chest ✗ | storage of urine ✗ | inflammation … ✗ |
| QCR: | no ✗ | intestine ✗ | treat brain diseases … ✗ |
| PubMedCLIP: | yes ✗ | spinal cord ✗ | nodule ✗ |

**Figure 5:** Examples from VQA-RAD where all models fail

For further analysis, we provide examples in Figure 5 from the VQA-RAD dataset, where all three models fail to yield the correct answer. We again observe that MEVF provides irrelevant answers about body organs that are not present in the image. QCR shows the same behaviour for the right-most example. For the image on the left, QCR and PubMedCLIP miscomprehend the question as a yes/no question. In spite of this, the fact that PubMedCLIP answers with "yes" illustrates that it has at least detected the "one" metastatic focus in the image. In comparison, QCR answers with "no", showing its troubles in interpreting the image and recognizing the metastatic focus. In the center example, answers provided by QCR and PubMedCLIP both appear to be relevant to the content of the given image. This suggests that the models have trouble understanding the semantics of the expression "periphery of the image" in the question. By invoking techniques such as Grad-CAM Selvaraju et al. [2017], we may be able to better understand what part of the image the model was focusing on before the classification layer. Finally, in the right-most example, QCR appears to misinterpret the content of the chest x-ray image and give suggestions for treating the brain. However, PubMedCLIP's answer, lung "nodule", seems to be at least relevant to the image, although it shows that it is having difficulties inferring the semantics of the question. Our observations reveal that these models still have shortcomings in understanding questions and correctly relating them to the images.

### 4.1 CLIP in MedVQA versus general-domain VQA

Using CLIP in general-domain VQA, as investigated in prior work Shen et al. [2021], evinces its effectiveness in comparison to previous ResNet-based encoders. In such settings, it has been observed that CLIP with a ResNet visual encoder outperforms using Vision Transformers. The authors hypothesize that this is due to Vision Transformers' weakness in visual localization. Furthermore, their reports show that larger back-end visual encoders in CLIP such as ResNet-101 and ResNet-50x4 result in bigger gains in accuracy.

In the MedVQA domain, we as well observe that PubMedCLIP outperforms the previous MAML-based visual encoders. However, using the bigger ResNet-50x4 model as the visual encoder in PubMedCLIP appears to lead to overfitting on medical images and it therefore performs slightly worse than the smaller version ResNet-50. Moreover, our experiments on the SLAKE dataset show that the Vision Transformer encoder in PubMedCLIP slightly outperforms using a ResNet-50. We showed that this gap stems mainly from differences in the underlying VQA data distributions. In the SLAKE dataset, the majority of the questions target a holistic overview of the image. In contrast, questions in VQA-RAD are primarily about the presence of an abnormality and require visual localization.

## 5 CONCLUSION

This work introduces PubMedCLIP as a pre-trained visual encoder for medical images and illustrates its effectiveness for the task of medical visual question answering. PubMedCLIP is trained using the image–caption pairs from thousands of PubMed articles. Our experiments on two benchmark MedVQA datasets demonstrate that PubMedCLIP outperforms previously used pre-trained visual encoders in MedVQA by up to 3%. Furthermore, our results reveal differences of underlying data distributions in the two benchmark datasets. We hope that our findings encourage future research to make real-world clinical image–text pairs publicly available for better development of vision–language representation

learning with cross-modal supervision in medical domain. In terms of future work, further analysis of these models using explainable AI techniques such as Grad-CAM visualizations can enable us to assess their regions of focus within the image from the class activation maps. Moreover, by releasing PubMedCLIP, we hope to enable research investigating to what extent it may contribute to additional medical use-cases such as image classification for medical diagnosis and radiology report generation.

## 6 ACKNOWLEDGMENT

We would like to thank Matthias Steinbrecher for his helpful comments and discussions.

## 7 CONFLICT OF INTEREST STATEMENT

None.

## References

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. VQA: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015.

Asma Ben Abacha, Vivek V Datla, Sadid A Hasan, Dina Demner-Fushman, and Henning Müller. Overview of the VQA-Med task at ImageCLEF 2020: Visual question answering and generation in the medical domain. In *CLEF (Working Notes)*, 2020.

Asma Ben Abacha, Mourad Sarrouti, Dina Demner-Fushman, Sadid A. Hasan, and Henning Müller. Overview of the VQA-Med task at ImageCLEF 2021: Visual question answering and generation in the medical domain. In *CLEF 2021 Working Notes*, CEUR Workshop Proceedings, Bucharest, Romania, September 21-24 2021. CEUR-WS.org.

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.

Haifan Gong, Ricong Huang, Guanqi Chen, and Guanbin Li. SYSU-HCP at VQA-Med 2021: A data-centric model with efficient training methodology for medical visual question answering. In *CLEF 2021 Working Notes*, volume 201. CEUR-WS, 2021.

Sedigheh Eslami, Gerard de Melo, and Christoph Meinel. TeamS at VQA-Med 2021: BBN-Orchestra for long-tailed medical visual question answering. In *Working Notes of CLEF 2021*, number 2936 in CEUR Workshop Proceedings, pages 1211–1217. CEUR-WS, September 21–24 2021. URL http://ceur-ws.org/Vol-2936/#paper-98.

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, 2021a.

Minh H Vu, Tommy Löfstedt, Tufve Nyholm, and Raphael Sznitman. A question-centric model for visual question answering in medical imaging. *IEEE transactions on Medical Imaging*, 39(9):2856–2868, 2020.

Li-Ming Zhan, Bo Liu, Lu Fan, Jiaxin Chen, and Xiao-Ming Wu. Medical visual question answering via conditional reasoning. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 2345–2354, 2020.

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.

Haiwei Pan, Shuning He, Kejia Zhang, Bo Qu, Chunling Chen, and Kun Shi. MuVAM: A multi-view attention-based model for medical visual question answering. *arXiv preprint arXiv:2107.03216*, 2021.

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, 2021b.

Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang. Bilinear attention networks. *arXiv preprint arXiv:1805.07932*, 2018.

Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alex Smola. Stacked attention networks for image question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 21–29, 2016.

Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*, pages 1126–1135. PMLR, 2017.

Jonathan Masci, Ueli Meier, Dan Cireşan, and Jürgen Schmidhuber. Stacked convolutional auto-encoders for hierarchical feature extraction. In *International conference on artificial neural networks*, pages 52–59. Springer, 2011.

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. *arXiv preprint arXiv:2103.00020*, 2021.

Jaemin Cho, Jie Lei, Hao Tan, and Mohit Bansal. Unifying vision-and-language tasks via text generation. *arXiv preprint arXiv:2102.02779*, 2021.

Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. Vl-bert: Pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.08530*, 2019.

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.

Sheng Shen, Liunian Harold Li, Hao Tan, Mohit Bansal, Anna Rohrbach, Kai-Wei Chang, Zhewei Yao, and Kurt Keutzer. How much can clip benefit vision-and-language tasks? *arXiv preprint arXiv:2107.06383*, 2021.

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.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=YicbFdNTTy.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. doi: 10.1109/CVPR.2016.90.

Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.

Damien Teney, Peter Anderson, Xiaodong He, and Anton Van Den Hengel. Tips and tricks for visual question answering: Learnings from the 2017 challenge. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4223–4232, 2018.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.

Hung Viet Pham, Shangshu Qian, Jiannan Wang, Thibaud Lutellier, Jonathan Rosenthal, Lin Tan, Yaoliang Yu, and Nachiappan Nagappan. Problems and opportunities in training deep learning software systems: an analysis of variance. In *Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering*, pages 771–783, 2020.

Qihang Yu, Yingda Xia, Yutong Bai, Yongyi Lu, Alan Yuille, and Wei Shen. Glance-and-gaze vision transformer. *arXiv preprint arXiv:2106.02277*, 2021.

Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-CAM: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.