

MMBERT: MULTIMODAL BERT PRETRAINING FOR IMPROVED MEDICAL VQA

Yash Khare^{*†} Viraj Bagal^{*‡} Minesh Mathew[†]
Adithi Devi^{††} U Deva Priyakumar[†] CV Jawahar[†]

[†] IIIT Hyderabad, India [‡]IISER Pune, India ^{††}Osmania Medical College, India

ABSTRACT

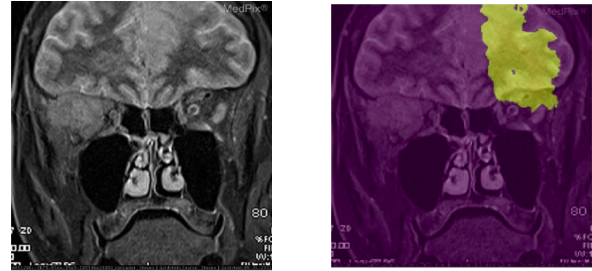
Images in the medical domain are fundamentally different from the general domain images. Consequently, it is infeasible to directly employ general domain Visual Question Answering (VQA) models for the medical domain. Additionally, medical images annotation is a costly and time-consuming process. To overcome these limitations, we propose a solution inspired by self-supervised pretraining of Transformer-style architectures for NLP, Vision and Language tasks. Our method involves learning richer medical image and text semantic representations using Masked Language Modeling (MLM) with image features as the pretext task on a large medical image+caption dataset. The proposed solution achieves new state-of-the-art performance on two VQA datasets for radiology images – VQA-Med 2019 and VQA-RAD, outperforming even the ensemble models of previous best solutions. Moreover, our solution provides attention maps which help in model interpretability.

Index Terms— medical VQA, multimodal BERT, vision and language

1. INTRODUCTION AND RELATED WORK

Visual question answering (VQA) on medical images aspires to build models that can answer diagnostically relevant natural language questions asked on medical images. It can provide valuable additional insights to medical professionals and can help the patients in the interpretation of their medical images. However, supervised learning algorithms require large labeled datasets for effective performance and a major drawback of VQA in the medical domain is the small size of existing datasets [1, 2, 3]. Since the annotations on medical images require the help of an expert, it is difficult to crowd-source and annotation cost is high. This motivates the usage of self-supervised pretraining methods.

Self-supervised pretraining of BERT-like architectures for proxy tasks like masked language modeling (masked LM) has been proven quite effective in Natural Language Processing (NLP) [4], Vision and Language [5, 6] space. The solution



Question: What imaging modality was used?
Answer: MR-T2 Weighted

Fig. 1: Example illustrating the attention map from our MMBERT model. For the given question, the model attends to grey matter, white matter and cerebrospinal fluid (CSF) and predicts the correct answer – 'MR-T2 Weighted'.

we propose - a Multimodal Medical BERT (MMBERT) is inspired by these approaches in the Vision and Language space. We first pretrain our MMBERT model on a set of medical images and their corresponding captions for the masked LM task. Later this model is finetuned for the VQA task.

Existing models for medical VQA are mostly inspired from the models developed for general domain VQA. Yan et al. [7] who are the winners of the VQA-Med 2019 challenge, use a Convolutional Neural Network (CNN) and BERT to extract image and question features respectively, followed by co-attention to fuse these features and a decoder to predict the answers. Ren et al. [8] propose a model called CGMVQA that uses a multimodal transformer architecture, similar to the proposed MMBERT. Zhan et al.[9] use a conditional reasoning framework for medical VQA on VQA-RAD dataset and they train a model separately for both the open-ended and closed-ended questions in the dataset.

Although the aforementioned methods obtain effective results, they do not use existing large multimodal medical datasets to learn better image and text representations. Our approach takes this into account and achieves new state-of-the-art performance on two medical VQA datasets. Our MMBERT, even with a single model for both the type of questions, yields better results than all the previous models on VQA-RAD. It also achieves a 5% improvement in Accuracy

^{*} Equal Contribution. Order decided by coin toss. Work done during an internship at IIIT Hyderabad.

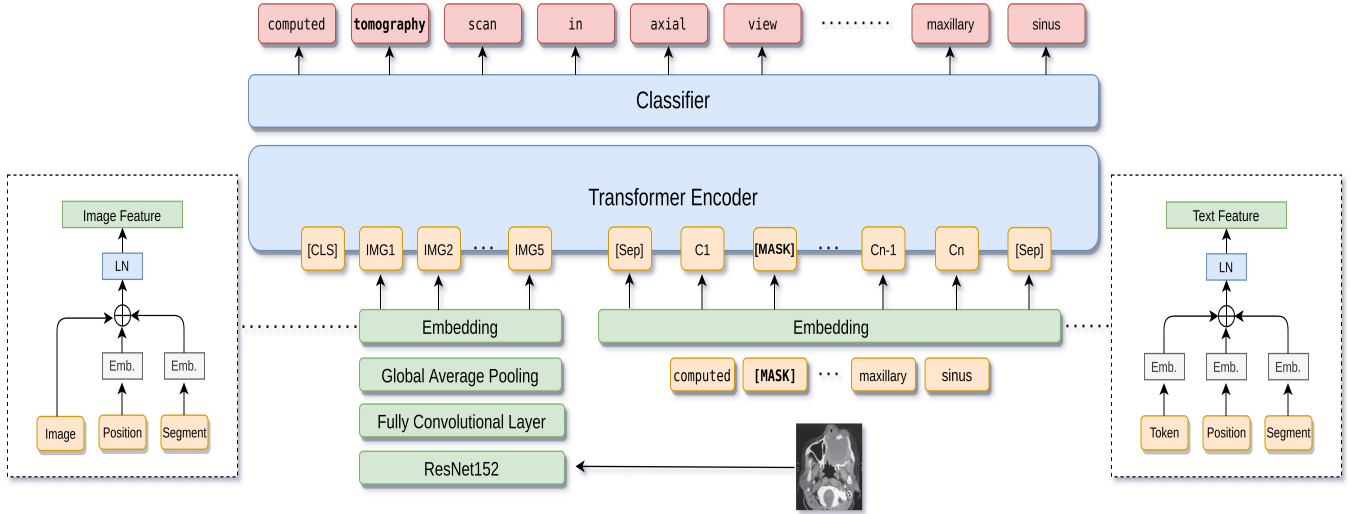


Fig. 2: Model architecture for MLM Pretraining task on image caption data. The image features are extracted from ResNet152 and passed through an embedding layer. The caption is tokenized and the keywords are masked with [Mask] tokens. The text embeddings are obtained by combining input, position and segment embeddings. The final embedding is passed through a transformer encoder. The encoder outputs are then passed to a classifier which predicts the masked words.

over the previous state-of-the-art model on VQA-Med 2019 dataset. Moreover, our model provides attention maps and as shown in Fig 1, it focuses on correct region (grey and white matter difference) to predict the modality of the image.

2. METHOD

Transfer learning is quite popular in machine learning. However, a shift in image data distribution might result in sub-optimal performance when using pretrained weights from general domain. Moreover, there are changes in co-occurrences of words in the medical text compared to the general domain text. These factors motivate the need for learning semantic representations of medical images and texts from scratch. Owing to the attention operation, recent Vision+Language architectures employ Transformer as the base architecture for learning effective representations and we do the same in our study.

2.1. Self-Attention

Self-Attention is the major mathematical operation occurring in the Transformer Encoder of Fig. 2. It allows attention to intra-modality and inter-modality features, i.e, image and text modality in our study, thus enhancing the semantics of the intermediate representations. Self-attention involves mapping a query vector to the weighted addition of the value vectors where the weights are obtained by scaling the dot product of the query and the key vectors. This is called 'Scaled Dot Product Attention' [10]. The query, key and value vectors corresponding to each token are computed by W_q , W_k and

W_v weight matrices and are represented together in the matrices Q , K and V respectively. The dot product of Q and K is scaled inversely by $\sqrt{d_k}$, where d_k is the dimension of query and key vectors.

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Instead of performing a single self-attention, the Transformer Encoder performs multiple self-attentions (multi-head attention) in parallel and concatenates the output. Multi-head attention provides better representations by attending to different representation subspaces at different positions.

2.2. Pretraining

A schematic of the MMBERT pretraining is shown in Figure 2. For image features, similar to the CGMVQA [8] we use ResNet152 [12] and extract features from different convolution layers. This helps in retaining information from different resolutions. We use BERT wordPiece tokenizer [4] for text tokenization. The sequence of 5 image features and the caption token embeddings together are provided as input to the BERT-like model. Unlike BERT_{BASE} [4] our model has only 4 BERT Layers and a total of 12 attention heads.

We use masked language modeling with image features as the pretraining task. In masked language modeling with image features, the task is to predict the original token in place of a [MASK] token with the usage of not only the accompanying text but also the image features. To ensure that the model learns to predict medical words from the context, we mask

Method	Dedicated Models	Modality		Plane		Organ		Abnormality		Yes/No		Overall	
		Acc.	BLEU	Acc.	BLEU	Acc.	BLEU	Acc.	BLEU	Acc.	BLEU	Acc.	BLEU
VGG16+BERT [7]	-	-	-	-	-	-	-	-	-	-	-	62.4	64.4
CGMVQA [8]	✓	80.5	85.6	80.8	81.3	72.8	76.9	1.7	1.7	75.0	75.0	60.0	61.9
CGMVQA Ens. [8]	✓	81.9	88.0	86.4	86.4	78.4	79.7	4.40	7.60	78.1	78.1	64.0	65.9
MMBERT General	✗	77.7	81.8	82.4	82.9	73.6	76.6	5.20	6.70	85.9	85.9	62.4	64.2
MMBERT NP	✓	80.6	85.6	81.6	82.1	71.2	74.4	4.30	5.70	78.1	78.1	60.2	62.7
MMBERT Exclusive	✓	83.3	86.2	86.4	86.4	76.8	80.7	14.0	16.0	87.5	87.5	67.2	69.0

Table 1: Results on VQA-Med 2019 dataset. Our method outperforms all previous methods that include methods with ensemble models in overall Accuracy and BLEU score. NP and Ens. refer to non-pretrained and ensemble models respectively.

Method	Dedicated Models	Accuracy		
		Open	Closed	Overall
MEVF+SAN [11]	-	40.7	74.1	60.8
MEVF+BAN [11]	-	43.9	75.1	62.7
CR [9]	✓	60.0	79.3	71.6
MMBERT General	✗	63.1	77.9	72.0

Table 2: Results on VQA-RAD dataset. Our method with single model for both open-ended and closed-ended question types outperforms all previous methods including methods with dedicated models for each question type in overall Accuracy.

only medical keywords (provided with the dataset) from the captions and leave the common words untouched.

2.3. Finetuning

We load the model with weights from pretraining and finetune it further on the train split of the respective medical VQA dataset. Instead of using [CLS] token representation from the last layer of the Transformer, we average the representation of each token obtained from the last layer and further pass it through dense layers for classification.

3. EXPERIMENTS AND RESULTS

3.1. Data

ROCO dataset contains over 81,000 radiology images with several medical imaging modalities. For pretraining, we use all the images, their corresponding captions and use the keywords for masking. VQA-Med 2019 [2] is a challenge dataset introduced as part of the ImageCLEF-VQA Med 2019 challenge. It contains radiology images and has four main categories of questions: Modality, Plane, Organ system and Abnormality. All the samples having Yes/No as the ground truth are considered as Yes/No category. The dataset includes a training set of 3200 medical images with 12,792 Question-Answer (QA) pairs, a validation set of 500 medical images with 2000 QA pairs and a test set of 500 medical images with 500 QA pairs. VQA-RAD has 315 images and 3515 questions of 11 types. 58% of questions are close-ended while the rest are open-ended.

3.2. Experiments

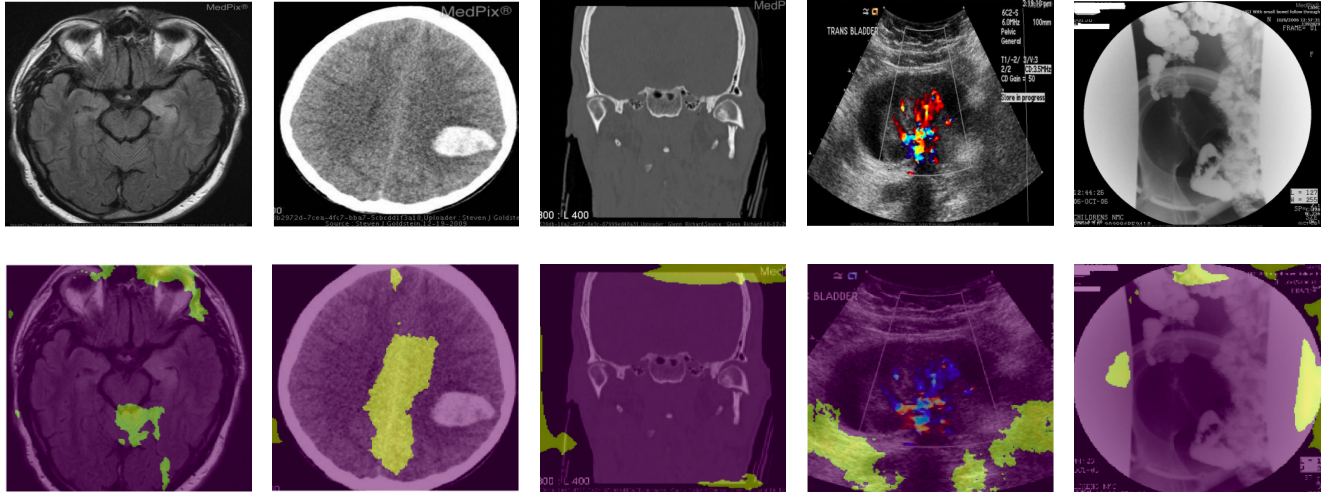
In our study, we primarily experiment with three different settings for the MMBERT: (i) MMBERT General: a model pretrained on ROCO and finetuned on all samples in the train split of the respective VQA dataset (ii) MMBERT Exclusive: an initial model pretrained on ROCO, which is further finetuned separately for different question categories. For example in case of VQA-Med 2019, we learn 5 different models, one for each question category and (iii) MMBERT Non-Pretrained (NP): Dedicated models for each question category but without pretraining on the ROCO dataset. At the time of the inference, for settings where there are dedicated models for each question category we first predict the question category using a $BERT_{BASE}$ classifier.

We train the models on a single NVIDIA RTX 2080Ti GPU. For pretraining and finetuning, we resize all images to 224×224 . We use image crop, rotation and color jitter for augmentation. For pretraining, we optimize the loss using Adam optimizer [13] with learning rate $2e - 5$ and reduce the learning rate by a factor of 0.1 if the validation loss does not improve for 5 consecutive epochs. For finetuning, we use the Adam optimizer, learning rate of $1e - 4$ and reduce the learning rate by a factor of 0.1 if validation loss does not improve for 10 consecutive epochs.

3.3. Results and Analysis

We use Accuracy and BLEU score to evaluate the VQA performance. Table 1 reports results on the VQA-Med 2019 dataset. Our MMBERT Exclusive achieves state-of-the-art results on the overall Accuracy and BLEU score, even surpassing CGMVQA Ens. which is an ensemble of 3 dedicated models for each category. Even our MMBERT General performs better than the CGMVQA Ens. on the Abnormality and Yes/No categories. Additionally, our MMBERT General outperforms single dedicated CGMVQA models in all categories but Modality.

In the Organ category, MMBERT Exclusive outperforms CGMVQA Ens. in BLEU but not in Accuracy. BLEU score is calculated by counting matching 1-gram in the predicted answer to the 1-gram in the ground truth. The comparison is made regardless of the order. This suggests that even though our model couldn't predict perfectly right answers, it could



<p>C: Organ Q: What organ is the image of? GT: skull and contents ME: skull and contents</p>	<p>C: Plane Q: What is the plane of the image? GT: axial ME: axial</p>	<p>C: Yes/No Q: Is this an MRI image? GT: no ME: no</p>	<p>C: Modality Q: What imaging method was used? GT: us-d - doppler ultrasound ME: us - ultrasound</p>	<p>C: Abnormality Q: What is abnormal in the image? GT: crohn's disease ME: fluoroscopic evaluation of small bowel in crohn's ileitis</p>
---	---	--	--	--

Fig. 3: ME refers to MMBERT Exclusive. The bottom row comprises attention maps for the corresponding top row images. In the Organ and Yes/No category, the model rightly attends to the bony part and soft tissue content to predict the right answer. In the Plane category, the model attends to the longitudinal fissure that is the key visual cue of the axial plane. The model fails to attend to the visual cue of doppler effect (colorful regions) in the Modality category. The Abnormality model surprisingly predicts a better answer than the ground truth by simultaneously predicting the modality, organ and abnormality.

predict more answers close to the ground truth than the CG-MVQA Ens. We find the opposite behaviour in the Modality category. When compared to MMBERT NP, we find that the pretraining increases the Accuracy and BLEU score by 7.2 and 9 points respectively.

Table 2 reports results on the VQA-RAD dataset. MMBERT General, which is a single model for both the question types in the dataset, outperforms the existing approaches including the ones which have a dedicated model for each question type.

3.4. Qualitative Analysis

Fig. 3 shows the category-wise qualitative results from MMBERT Exclusive. The top row comprises the original images while the bottom row comprises the attention maps obtained from our model. The attention maps highlight the regions in the image which contribute the most to the prediction. In the Organ and Yes/No category, the model rightly attends to the skull (the bony part) and its contents (brain tissues) to predict the right answer. In the Plane category, the model attends to the longitudinal fissure which is the key visual cue in identifying the axial plane as it separates the brain into Right and Left hemispheres. In the Modality category, the model attends to the soft tissue and fluid part of the image and is able to cor-

rectly predict that it is an ultrasound image. However, it fails to attend to the visual cue of the Doppler (the colour region) and hence fails to correctly answer. Surprisingly, in the case of Abnormality category, our model predicts a better answer than the ground truth. Here, it is simultaneously predicting the modality (fluoroscopy), organ (bowel) and the abnormality (Crohn’s ileitis).

Medical experts find it difficult to make a correct diagnosis of abnormalities from a single image. They often resort to multiple sections (slices), planes, and other evidences. On closely analyzing our results we see that our model predicts abnormalities which could have also been a differential diagnosis for a human expert. However, our quantitative evaluation protocol does not take this into consideration.

4. CONCLUSION

In this work, we propose to pretrain Multimodal Medical BERT (MMBERT) on ROCO dataset with masked language modeling using image features for medical VQA. We finetune it on VQA-RAD and VQA-Med 2019 datasets and achieve new state-of-the-art results on these datasets. Moreover, qualitative results show that our models can rightly attend to the image regions for prediction.

5. COMPLIANCE WITH ETHICAL STANDARDS

Datasets used in the research did not require ethical approval.

6. ACKNOWLEDGMENTS

We have no financial or non-financial interests to disclose.

7. REFERENCES

- [1] 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*, 2018.
- [2] Asma Ben Abacha, Sadid A. Hasan, Vivek V. 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 Working Notes*, 2019, CEUR Workshop Proceedings.
- [3] Xuehai He, Yichen Zhang, Luntian Mou, Eric Xing, and Pengtao Xie, "PathVQA: 30000+ questions for medical visual question answering," *arXiv preprint arXiv:2003.10286*, 2020.
- [4] 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.
- [5] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee, "Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks," in *NeurIPS*, 2019.
- [6] Chen Yen-Chun, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu, "Uniter: Universal image-text representation learning," in *ECCV*, 2020.
- [7] Xin Yan, Lin Li, Chulin Xie, Jun Xiao, and Lin Gu, "Zhejiang university at ImageCLEF 2019 visual question answering in the medical domain.," in *CLEF (Working Notes)*, 2019.
- [8] Fuji Ren and Yangyang Zhou, "CGMVQA: A New Classification and Generative Model for Medical Visual Question Answering," in *IEEE Access*, 2020.
- [9] Li-Ming Zhan, Bo Liu, Lu Fan, Jiabin Chen, and Xiao-Ming Wu, "Medical visual question answering via conditional reasoning," in *ACM*, 2020.
- [10] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin, "Attention is all you need," 2017.
- [11] 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," *MICCAI*, 2019.
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," *arXiv preprint arXiv:1512.03385*, 2015.
- [13] Diederik P Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.