# Cross-Modal Self-Attention with Multi-Task Pre-Training for Medical Visual Question Answering

Haifan Gong*
Sun Yat-sen University
gonghf@mail2.sysu.edu.cn

Guanqi Chen*
Sun Yat-sen University
chengq26@mail2.sysu.edu.cn

Sishuo Liu
The University of Hong Kong
sishuo@hku.hk

Yizhou Yu
The University of Hong Kong
yizhouy@acm.org

Guanbin Li†
Sun Yat-sen University
liguanbin@mail.sysu.edu.cn

## ABSTRACT

Due to the severe lack of labeled data, existing methods of medical visual question answering usually rely on transfer learning to obtain effective image feature representation and use cross-modal fusion of visual and linguistic features to achieve question-related answer prediction. These two phases are performed independently and without considering the compatibility and applicability of the pre-trained features for cross-modal fusion. Thus, we reformulate image feature pre-training as a multi-task learning paradigm and witness its extraordinary superiority, forcing it to take into account the applicability of features for the specific image comprehension task. Furthermore, we introduce a cross-modal self-attention (CMSA) module to selectively capture the long-range contextual relevance for more effective fusion of visual and linguistic features. Experimental results demonstrate that the proposed method outperforms existing state-of-the-art methods. Our code and models are available at https://github.com/haifangong/CMSA-MTPT-4-MedicalVQA.

## CCS CONCEPTS

• **Computing methodologies** → **Computer vision**; **Natural language processing**.

## KEYWORDS

Visual question answering, transfer learning, multi-task learning, self-attention

*Both authors contributed equally to this work. Guanbin Li is the corresponding author.

## 1 INTRODUCTION

The medical visual question answering (VQA) aims to answer the questions of images in the medical domain. The common setting of medical VQA is to retrieval the answer from the answer set which best fits the given question and the image. With the expectation that VQA systems can not only provide clinicians with clinical decision support, but also help patients better understand their conditions based on medical images, several medical VQA datasets[1, 2, 14] have been proposed. Since the questions and related answers are automatically generated, the medical VQA datasets produced by [1, 2] inevitably contain noise, which may not be the optimal choice for clinical decision support. Different from [1, 2], the VQA-RAD dataset[14] was manually constructed where clinicians naturally create the question-answer pairs about the radiology images. Unfortunately, there are only 315 images in the VQA-RAD dataset, while the out-performing deep learning based algorithms require large-scale data for training. Therefore, many works[3, 14, 17, 25] adopt transfer learning to solve this problem, relying on external data for pre-training of image features before the training of feature fusion and answer prediction. However, those works neglect considerations for the compatibility and applicability of the pre-trained features for cross-modal fusion. In addition, compared to the VQA task of natural images, there are several unique challenges in the medical VQA domain, including semantic parsing of medical terminology, more complex cross-modal semantic alignment and fusion due to low contrast of medical images, and the multi-modal characteristics of medical images (i.e., CT, MRI, X-Ray).

Based on the above concerns, we propose to reformulate image feature pre-training as a multi-task learning paradigm, forcing it to take into account the applicability of features for both the specific image comprehension task and our proposed cross-modal fusion module. This has been proven to make more effective use of external data to better overcome the problem of data scarcity in medical VQA. Secondly, a tailor-designed cross-modal self-attention (CMSA) module is used to effectively fuse the visual and linguistic features by learning and leveraging their long-range contextual relevance, which effectively compensates for the low contrast and weak local feature representation in medical images through contextual information enhancement and complementation. Last but not least, we achieve the state-of-the-art performance on the VQA-RAD dataset.
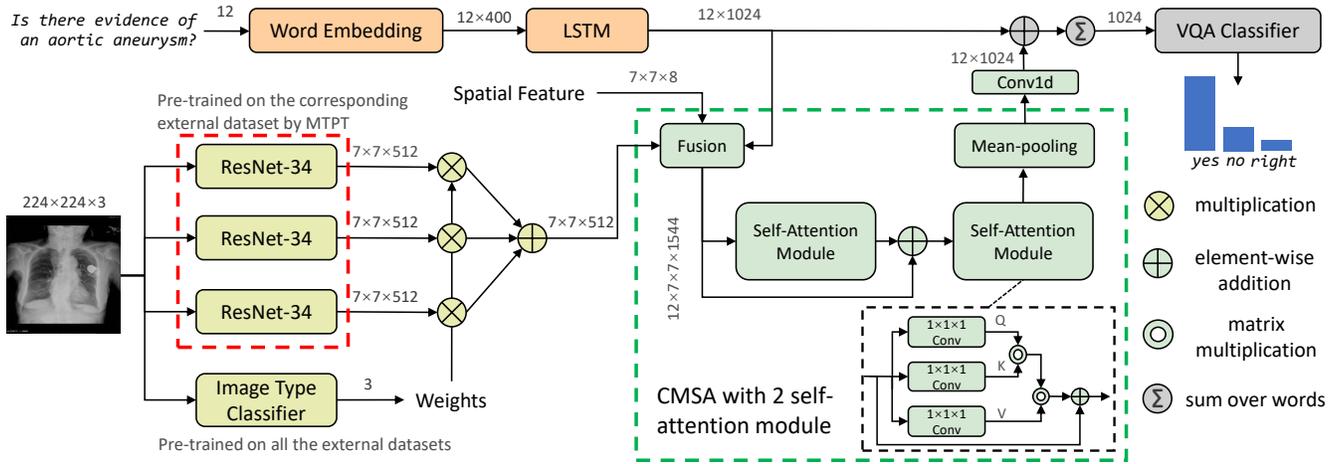
**Figure 1: Overview of the proposed medical VQA model. Our method consists of four components (with different colors in the figure): image feature extractor, question encoder, cross-modal self-attention (CMSA) module, and answer predictor.**

## 2 RELATED WORK

### 2.1 Visual Question Answering

With the prosperity of deep learning, VQA has received extensive attention in recent years and has made great progress, mainly benefiting from the strong bottom-up feature representation based on deep CNNs[4] and the cross-modal feature alignment and fusion techniques[4, 8, 13, 26, 28]. Anderson et al. [4] proposed a bottom-up mechanism implemented by Faster R-CNN[19] to extract object-level representation for the input image, which achieved great success in both VQA and image captioning. From the perspective of cross-modal feature fusion, methods can be roughly divided into two main categories, including the attention based methods and multi-modal joint embedding. Anderson et al.[4] and Yang et al.[26] developed different attention modules to adaptively attend on the relevant image regions based on the question representation. Kim et al. [8, 13, 28] proposed to employ the compact bilinear pooling methods to combine the visual and linguistic features.

For medical VQA, the current common methods [12, 14, 17, 18, 22, 23, 25] are to use CNN for image feature representation and leverage LSTM [10] or transformer-based methods (e.g., Bert [7], BioBert [15]) to extract features for the given question. Varieties of general cross-modal fusion strategies (e.g., SAN[26], BAN[13], and MFB[28]) are applied for feature fusion followed by the ultimate answer prediction. Compared to general VQA, medical VQA systems are required to comprehend medical terminology and focus on the corresponding visual content in the image. However, existing medical VQA methods do not realize the significance of these problems and directly borrow the general VQA technologies, which caused bottlenecks in the prediction accuracy of the models.

### 2.2 Transfer Learning

Due to the limitation of medical VQA data, many works rely on transfer learning to obtain effective image feature representation. In [5, 14, 22, 25], they use a CNN which is pre-trained on ImageNet[20] to encode medical image, such as VGGNet and ResNet. Allaouzi *et*

*al.* [3] utilizes CheXpert [11], a large dataset of chest radiographs, to pre-train a DenseNet-121 as the visual feature encoder. Nguyen *et al.* [17] leverages a large scale of unlabeled medical images to pre-train its unsupervised denoising auto-encoder via a reconstruction task. However, there are no existing works attempt to consider the compatibility and applicability of the pre-trained features for cross-modal fusion, which is the emphasis of VQA models.

## 3 METHODOLOGY

The proposed medical VQA framework is shown in Figure 1, which includes a multi-task pre-training paradigm for more effective medical image representation learning, a cross-modal self-attention module for feature fusion, and the ultimate VQA classifier for question-related answer prediction. In this section, we elaborate the proposed multi-task pre-training method and the medical VQA model.

### 3.1 Multi-Task Pre-Training

During multi-task pre-training, our model is jointly trained with two separate tasks, including a regular image understanding task and a tailor-designed task for question-image compatibility test. The latter is defined as a binary classification task, which requires the model to determine whether the question is related to and suitable for a given image. For example, the question "Are the lungs normal size?" is suitable for a chest image rather than an abdominal image. For a given image from the external dataset, we randomly select a question from the VQA-RAD to form a question-image pair. The label for question-image compatibility testing is constructed by querying whether there exists a pair of the selected question and an image whose type is the same as the given image in VQA-RAD.

As shown in Figure 2, during pre-training with external datasets, we use ResNet-34 [9] as a backbone to capture the visual feature of the input image, and a decoder of symmetric structure for segmentation and a 3-layer MLP for image classification. For question-image compatibility testing, we use the proposed cross-modal self-attention (CMSA) module for feature fusion, which will be detailed
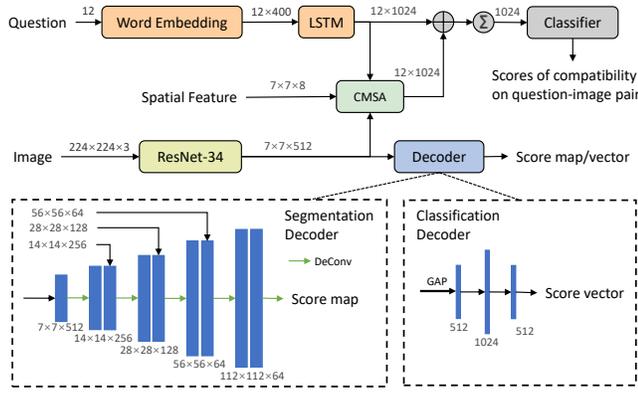
**Figure 2: Multi-Task Pre-Training: the model is jointly trained with an image understanding task and a question-image compatibility task. Depending on the dataset-specific image understanding task, the decoder can be selected as a fully convolutional network or a fully connected network.**

later. It is worth noting that the CMSA module used here only contains one self-attention module without repeating, because we want this pre-training task focus on the representation learning of the image encoders rather than feature fusion. Formally, the multi-task loss function is defined as:

$$L = L_{spe} + L_{com} \tag{1}$$

where $L_{spe}$ and $L_{com}$ are the cross-entropy loss for the specific image understanding task and the question-image compatibility task, respectively.

### 3.2 Our Medical VQA Model

Our medical VQA model consists of four parts: image encoding for capturing visual features of the given medical image, question encoding for extracting language features of the given question, cross-modal self-attention module for visual-language feature fusion, and answer prediction. We use a multi-task loss $L$ to train the proposed medical VQA model in an end-to-end fashion:

$$L = L_{vqa} + \alpha L_{type} \tag{2}$$

where $L_{vqa}$ and $L_{type}$ are the cross-entropy loss for classification based answer prediction and image type classification. $\alpha$ is a hyperparameter for balancing the two loss terms, which is set to 0.5.

*3.2.1 Image encoding:* Since clinicians use different imaging techniques to precisely diagnose the diseases of different organs, we use three separate ResNet-34 networks pre-trained on the corresponding external datasets to capture visual features of MRI, CT, X-Ray images, respectively. Then we use a classifier to determine the type of medical image and select the corresponding visual feature in a soft manner:

$$v = w_1 v_a + w_2 v_h + w_3 v_c \tag{3}$$

where $v$ denotes the final visual feature, $v_a$, $v_h$ and $v_c$ denote the output features from the encoder corresponding to the abdomen, head and chest images respectively. $w$ is the output vector of image type classifier with $\sum_{l=1}^{3} w_l = 1$, which represents the corresponding

weight of each medical image type. Besides, to better comprehend and answer questions related to the localization of local images, we follow [27] to obtain an 8-D spatial feature map $s$ with the same resolution as the visual feature $v$. The spatial vector at each position in the spatial feature map $s$ encodes the normalized coordinates of top-left, center, bottom-right, width and height of the grid.

*3.2.2 Question encoding:* Following the previous work [17], each input question is trimmed to a maximum of 12 words, and it is zero-padded when its length is less than 12. Each word is represented as a concatenation of a 200-D BioWordVec [30] word embedding and another 200-D augmenting embedding from the VQA-RAD. BioWordVec is a pre-trained biomedical word embedding based on PubMed and MeSH. Each 400-D word embedding vector is further fed into a LSTM to obtain the question embedding $q \in \mathbb{R}^{12 \times 1024}$.

*3.2.3 Cross-Modal Self-Attention:* Before cross-modal fusion, we are given the visual feature $v \in \mathbb{R}^{7 \times 7 \times 512}$, spatial feature $s \in \mathbb{R}^{7 \times 7 \times 8}$ and question embedding $q \in \mathbb{R}^{12 \times 1024}$. For each word in the question, we concatenate its representation with visual and spatial features at each spatial location to produce a feature map $f \in \mathbb{R}^{7 \times 7 \times 1544}$. Then, we collect all the concatenated feature maps to obtain a multi-modal feature map $F \in \mathbb{R}^{12 \times 7 \times 7 \times 1544}$. Inspired by the effectiveness of self-attention[21, 24] in capturing non-local context, we design our cross-modal alignment and fusion method.

Firstly, we linearly transform the multi-modal feature map $F$ to produce three feature maps $Q, K, V \in \mathbb{R}^{12 \times 7 \times 7 \times 772}$ through three $1 \times 1 \times 1$ convolutional layers . We reshape them to the dimension $\mathbb{R}^{588 \times 772}$, and use $Q$ and $K$ to compute the attention map $A$:

$$A = softmax(QK^T) \tag{4}$$

where $A \in \mathbb{R}^{588 \times 588}$ indicates the correlation of features in different positions. We multiply the attention map $A$ and the feature map $V$ to obtain the enhanced multi-modal representation $F' \in \mathbb{R}^{588 \times 772}$:

$$F' = AV. \tag{5}$$

Next, we turn the dimension of $F'$ to $\mathbb{R}^{12 \times 7 \times 7 \times 1544}$ through reshaping and $1 \times 1 \times 1$ convolution. The above operations are shown in the Figure 1 named self-attention module. Inspired by the 'glimpse' in BAN[13], we repeat the self-attention module again with residual connection. The final multi-modal representation $\hat{F} \in \mathbb{R}^{12 \times 1544}$ is obtained by applying a mean-pooling operation to the output of the residual connection between $F'$ and $F$ over all spatial locations:

$$\hat{F}_i = \frac{\sum_{j=1}^{7} \sum_{k=1}^{7} \left( F'_{ijk} + F_{ijk} \right)}{7 \times 7} \tag{6}$$

where $i$, $j$, $k$ are the indices of the number of words, height and width of the feature map. Then $\hat{F}$ is transformed to the same dimension as the question embedding $q$ with a linear layer.

*3.2.4 Answer prediction:* The joint representation $\hat{F}$ is added element-wise with question embedding $q$, and it is summed over all words in the question. Finally, we feed it into a 2-layer MLP for answer prediction. Prediction score $s$ of the answer is calculated by:

$$s = MLP(\sum_{i=1}^{12} (\hat{F}_i + q_i)). \tag{7}$$

**Table 1: Comparisons with the state-of-the-art methods on the VQA-RAD test set.** $para^\star$ **means using not only the "freeform" but also the "para" answer type in the test set.**

| Methods | Open | Closed | All |
|---|---|---|---|
| SAN-RAD [10] | 24.2% | 57.2% | 44.0% |
| MCB-RAD [10] | 25.4% | 60.6% | 46.5% |
| SAN-MEVF [11] | 40.7% | 74.1% | 60.8% |
| BAN-MEVF [11] | 43.9% | 75.1% | 62.6% |
| Ours | 56.1% | 77.3% | 68.8% |
| BAN-CR-para$^\star$ [22] | 60.0% | 79.3% | 71.6% |
| Ours-para$^\star$ | 61.5% | 80.9% | 73.2% |

**Table 2: Ablation study on the 3 external validation set.**

| Methods | Abdominal CT | | Brain MRI | | Chest X-Ray | |
|---|---|---|---|---|---|---|
| | mIOU | $acc_{com}$ | $acc_{cls}$ | $acc_{com}$ | $acc_{cls}$ | $acc_{com}$ |
| Baseline | 0.682 | - | 98.4% | - | 97.8% | - |
| MTPT | 0.710 | 78.7% | 98.4% | 89.1% | 98.7% | 83.6% |

## 4 EXPERIMENTS

### 4.1 Datasets and Metrics

*4.1.1 Datasets:* The proposed method is evaluated on the VQA-RAD dataset[14], which contains 315 radiological images with 3064 training questions and 451 test questions. We resort to three external datasets to pre-train the visual encoders of different image types, including abdominal CT [1], brain MRI[6] and chest X-Ray [2]. The abdominal CT dataset includes 2178 images of 13 classes for multi-organ segmentation, where we use 2070 images for training and 108 images for validation. The brain MRI dataset comtains 3 types of brain tumors with 3604 images. We divide it into 3000 images and 64 images for training and validation, respectively. The chest X-Ray dataset contains 5232 images of 'pneumonia' or 'normal', which is split into 5000 images for training and 232 images for validation.

*4.1.2 Evaluation metrics:* We use accuracy as the metric for the VQA task and the classification task during pre-training, where $acc_{cls}$ and $acc_{com}$ refer to the accuracy of the image classification task and the question-image compatibility task respectively. Mean Intersection-over-Union (mIoU) is the criteria for segmentation.

### 4.2 Comparison with the state-of-the-art

As shown in Table 1, our proposed method is compared with 5 existing state-of-the-art approaches[14, 17, 29], and achieves the highest accuracy on both open-ended and closed-ended VQA. Compared with the advanced approach BAN-MEVF[17] using external datasets, the proposed method outperforms it by 12.2% and 2.2% w.r.t. accuracy on open-ended and closed-ended VQA, respectively. For the sake of fairness, the comparison between the proposed and the current best model BAN-CR[29] is based on the "freeform" and "para" questions as the setting in their code. The proposed method

**Table 3: Ablation study on the VQA-RAD test set.**

| Methods | Open | Closed | All |
|---|---|---|---|
| INPT-CMSA | 30.9% | 73.5% | 56.5% |
| STPT-CMSA | 41.5% | 74.1% | 61.0% |
| MTPT-BAN | 56.1% | 75.7% | 67.9% |
| MTPT-CMSA | 56.1% | 77.3% | 68.8% |

outperforms BAN-CR[29] by 1.6% w.r.t. accuracy on all the questions. Nevertheless, the proposed method could be combined with the conditional reasoning[29] to gain a further improvement.

### 4.3 Ablation Study

To explore the effectiveness of the multi-task pre-training method diagram, we compare it with a single-task pre-training method which only pre-train on the external datasets for the original image classification or segmentation task. In Table 2, 'Baseline' represents the single-task pre-training method. 'MTPT' denotes the designed method with BioWordVec word embedding. The results show that our proposed multi-task pre-training method can slightly improve the performance of each specific image understanding task.

After pre-training the visual encoders, we load the pre-trained weights to train the whole VQA model on VQA-RAD. To clearly illustrate the ablation study in Table 3, we give the definitions as: (1) 'INPT' uses three pre-trained ResNet-34 on ImageNet as visual encoders. 'STPT' initializes the visual encoders with single-task pre-training while 'MTPT' loads the weights of visual encoders from the multi-task pre-training. (2)'CMSA' uses the proposed 'CMSA' for feature fusion while 'BAN' applies BAN[13] for feature fusion.

From Table 3, 'MTPT-CMSA' outperforms 'STPT-CMSA' significantly with the same external datasets for pre-training, which suggests that the pre-trained visual features from our multi-task learning paradigm are more suitable for our CMSA module to obtain an effective multi-modal representation. Furthermore, the proposed 'CMSA' feature fusion surpasses the 'BAN' feature fusion method by capturing the long-range contextual relevance.

## 5 CONCLUSION

This paper introduces a distinguished medical VQA framework which is based on multi-task pre-training paradigm for more effective medical image representation learning. Moreover, the proposed CMSA module effectively fusion of visual and language features by capturing the long-range contextual relevance. Experimental results verify that our proposed method can leverage external data more effectively to overcome the limitation of medical VQA data. In the future, we gonna focus on integrating domain knowledge into medical VQA on the recent knowledge based dataset [16] for the interpretable medical application.

---

[1]https://www.synapse.org/#!Synapse:syn3193805/wiki/217753
[2]https://www.kaggle.com/paultimothymooney/chest-xray-pneumonia

# REFERENCES

[1] Asma Ben Abacha, Soumya Gayen, Jason J Lau, Sivaramakrishnan Rajaraman, and Dina Demner-Fushman. 2018. NLM at ImageCLEF 2018 Visual Question Answering in the Medical Domain.. In *CLEF (Working Notes)*.

[2] Asma Ben Abacha, Sadid A Hasan, Vivek V Datla, Joey Liu, Dina Demner-Fushman, and Henning Müller. 2019. VQA-Med: Overview of the medical visual question answering task at imageclef 2019. In *CLEF2019 Working Notes. CEUR Workshop Proceedings*. 09–12.

[3] Imane Allaouzi, B Benamrou, and MB Ahmed. 2019. An encoder-decoder model for visual question answering in the medical domain. *Working Notes of CLEF* (2019).

[4] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 6077–6086.

[5] Guanqi Chen, Haifan Gong, and Guanbin Li. 2020. HCP-MIC at VQA-Med 2020: Effective Visual Representation for Medical Visual Question Answering. In *Working Notes of CLEF 2020 - Conference and Labs of the Evaluation Forum, Thessaloniki, Greece, September 22-25, 2020 (CEUR Workshop Proceedings, Vol. 2696)*.

[6] Jun Cheng. 2017. brain tumor dataset. (4 2017). https://doi.org/10.6084/m9.figshare.1512427.v5

[7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. 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)*. 4171–4186.

[8] Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach. 2016. Multimodal Compact Bilinear Pooling for Visual Question Answering and Visual Grounding. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*. The Association for Computational Linguistics, 457–468.

[9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.

[10] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.

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

[12] Bumjun Jung, Lin Gu, and Tatsuya Harada. 2020. bumjun_jung at VQA-Med 2020: VQA Model Based on Feature Extraction and Multi-modal Feature Fusion. In *Working Notes of CLEF 2020 - Conference and Labs of the Evaluation Forum, Thessaloniki, Greece, September 22-25, 2020 (CEUR Workshop Proceedings, Vol. 2696)*. CEUR-WS.org.

[13] Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang. 2018. Bilinear Attention Networks. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*. 1571–1581.

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

[15] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics* 36, 4 (2020), 1234–1240.

[16] Bo Liu, Li-Ming Zhan, Li Xu, Lin Ma, Yan Yang, and Xiao-Ming Wu. 2021. SLAKE: A Semantically-Labeled Knowledge-Enhanced Dataset for Medical Visual Question Answering. In *18th IEEE International Symposium on Biomedical Imaging, ISBI 2021, Nice, Antipolis, France, April 13-16, 2021*. IEEE, pp.1–5.

[17] Binh D Nguyen, Thanh-Toan Do, Binh X Nguyen, Tuong Do, Erman Tjiputra, and Quang D Tran. 2019. Overcoming Data Limitation in Medical Visual Question Answering. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 522–530.

[18] Fuji Ren and Yangyang Zhou. 2020. CGMVQA: A New Classification and Generative Model for Medical Visual Question Answering. *IEEE Access* 8 (2020), 50626–50636.

[19] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. 2017. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. *IEEE Trans. Pattern Anal. Mach. Intell.* 39, 6 (2017), 1137–1149.

[20] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. 2015. Imagenet large scale visual recognition challenge. *International journal of computer vision* 115, 3 (2015), 211–252.

[21] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.

[22] Minh Vu, Raphael Sznitman, Tufve Nyholm, and Tommy Löfstedt. 2019. Ensemble of streamlined bilinear visual question answering models for the imageclef 2019 challenge in the medical domain. In *CLEF 2019*, Vol. 2380.

[23] Minh H Vu, Tommy Löfstedt, Tufve Nyholm, and Raphael Sznitman. 2020. A Question-Centric Model for Visual Question Answering in Medical Imaging. *IEEE transactions on medical imaging* 39, 9 (2020), 2856–2868.

[24] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. 2018. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7794–7803.

[25] Xin Yan, Lin Li, Chulin Xie, Jun Xiao, and Lin Gu. 2019. Zhejiang university at imageclef 2019 visual question answering in the medical domain. *Working Notes of CLEF* (2019).

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

[27] Linwei Ye, Mrigank Rochan, Zhi Liu, and Yang Wang. 2019. Cross-modal self-attention network for referring image segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 10502–10511.

[28] Zhou Yu, Jun Yu, Jianping Fan, and Dacheng Tao. 2017. Multi-modal factorized bilinear pooling with co-attention learning for visual question answering. In *Proceedings of the IEEE international conference on computer vision*. 1821–1830.

[29] Li-Ming Zhan, Bo Liu, Lu Fan, Jiaxin Chen, and Xiao-Ming Wu. 2020. Medical Visual Question Answering via Conditional Reasoning. In *Proceedings of the 28th ACM International Conference on Multimedia*. Association for Computing Machinery, New York, NY, USA, 2345–2354.

[30] Yijia Zhang, Qingyu Chen, Zhihao Yang, Hongfei Lin, and Zhiyong Lu. 2019. BioWordVec, improving biomedical word embeddings with subword information and MeSH. *Scientific data* 6, 1 (2019), 1–9.