

# Q2ATransformer: Improving Medical VQA via an Answer Querying Decoder

Yunyi Liu<sup>1</sup>, Zhanyu Wang<sup>1</sup>, Dong Xu<sup>2</sup>, and Luping Zhou<sup>1</sup>

<sup>1</sup> The University of Sydney, Sydney, NSW, Australia  
{yunyi.liu, zhanyu.wang, luping.zhou}@sydney.edu.au

<sup>2</sup> The University of Hong Kong, Hong Kong SAR  
dongxu@cs.hku.hk

**Abstract.** Medical Visual Question Answering (VQA) systems play a supporting role to understand clinic-relevant information carried by medical images. The questions to a medical image include two categories: close-end (such as Yes/No question) and open-end. To obtain answers, the majority of the existing medical VQA methods rely on classification approaches, while a few works attempt to use generation approaches or a mixture of the two to process the two kinds of questions separately (classification for the close-end and generation for the open-end). The classification approaches are relatively simple but perform poorly on long open-end questions, while the generation approaches face the challenge of generating many non-existent answers, resulting in low accuracy rates. To bridge this gap, in this paper, we propose a new Transformer based framework for medical VQA (named as Q2ATransformer), which integrates the advantages of both the classification and the generation approaches and provides a unified treatment for the close-end and open-end questions. Specifically, we introduce an additional Transformer decoder with a set of learnable candidate answer embeddings to query the existence of each answer class to a given image-question pair. Through the Transformer attention, the candidate answer embeddings interact with the fused features of the image-question pair to make the decision. In this way, despite being a classification-based approach, our method provides a mechanism to interact with the answer information for prediction like the generation-based approaches. On the other hand, by classification, we mitigate the task difficulty by reducing the search space of answers. Our method achieves new state-of-the-art performance on two medical VQA benchmarks. Especially, for the open-end questions, we achieve 79.19% on VQA-RAD and 54.85% on PathVQA, with 16.09% and 41.45% absolute improvements, respectively.

**Keywords:** Medical VQA · Attention Mechanism · Classification.

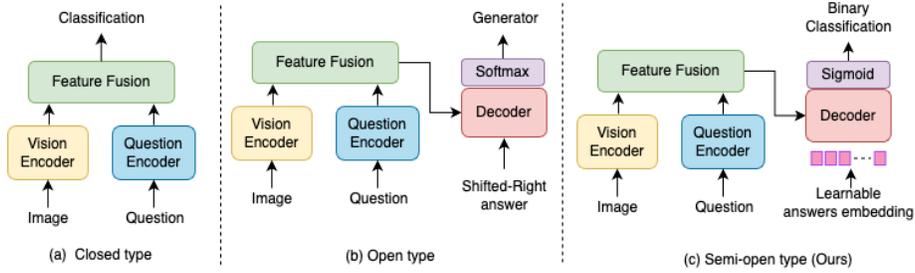
## 1 Introduction

Visual question answering (VQA) is known to be a challenging AI task that answers image-related questions based on image content. This process involves both

image and natural language processing techniques and usually comprises of four key components: extracting image features, extracting question features, integrating features, and answering. Recent years have witnessed significant progress in this field [8,15]. Medical VQA is a natural extension of VQA to medical images accompanied by clinic-relevant questions. Through questioning and answering, it offers a user-friendly way to assist clinic decisions. The questions in medical VQA could be either close-end, such as Yes/No questions, or open-end.

Medical VQA is still in its early stage of development and the current performance is far from being satisfying. Most existing methods [13,4,7,6,5] could be referred to as closed-type approaches, as illustrated in Figure 1 (a), which treat each answer as a class and apply a classification model directly to the fused features of the input image-question pair to predict answers. The advantage of such approaches is that by treating VQA as classification tasks, they reduce the complexity of the task and make the answer search space smaller. Despite the good performance on Yes/No questions, closed-type approaches are difficult to accurately predict the answer for open-end questions that are much longer and more varied than the close-end ones. On the other hand, a few works [2,9] treat VQA as a generation task and employ generation-based approaches to produce answers word by word. They are referred to as the open-type approaches in Figure 1(b). In these approaches, current word generation usually depends on previous words of the answer. Therefore, these approaches allow the image-question features to interact with the answer information for the prediction, potentially improving the long answer prediction. However, due to the tremendous search space of the generated answers, these approaches tend to produce many non-existent answers, leading to low accuracy rates, therefore are not currently the mainstream of medical VQA. Although there are attempts to combine these two types of approaches [14], they straightforwardly treat close-end and open-end questions separately, e.g., classification for close-end ones and generation for open-end ones.

To bridge the research gap and promote medical VQA, we introduce a new model framework Q2ATransformer and refer it to as semi-open type, as shown in Figure 1 (c). By semi-open, we keep adopting classification-based approaches to make the answer search space small, and at the same time introduce the learning of answer semantics so that the fused image-question features and the answer semantics could interact for better prediction, like the generation-based approaches. Our model mitigates the shortcomings of the classification-based closed-type framework while enjoying the advantages of the generation-based open-type framework. To achieve this, we introduce a set of learnable candidate answer embeddings and let the image-question feature interact with the candidate answer embeddings by sending them through a transformer decoder. In the decoder, the candidate answer embeddings work as a query to calculate their relationships with the fused image-question features to decide the existence of the answer classes. By this, our classification considers the interaction of answer information and the fused image-question features, which is different from the existing classification-based approaches. Compared with the generation-based



**Fig. 1.** Paradigms for medical VQA frameworks. (a) Closed-type framework treats VQA as predicting answer classes, where a classifier is built directly on top of the fused image-question features. (b) Open-type framework is generation-based, where the fused image-question features interact with the previous words of the answer to generate the next word of the answer through a text decoder. (c) Our proposed semi-open framework learns candidate answer embeddings through a decoder, where they interact with the fused image-question features to improve the prediction of answer classes.

open-type approaches, our model reduces the task difficulty and significantly improves the accuracy rates.

Last but not the least, our model provides a uniform treatment for both the close-end and the open-end questions.

The main contribution of this paper could be summarized as follows.

First, we proposed a framework of semi-open type for medical VQA, which bridges the advantages of both the classification-based closed-type framework and the generation-based open-type frameworks in medical VQA literature. This is achieved by a designed mechanism to learn and make use of candidate answer embedding through a transformer decoder while limiting the search scope of answers through classification.

Second, we proposed a Cross-modality Fusion Network (CMAN) to effectively fuse the image and question features. It directly concatenates the two modal features instead of conducting matrix multiplication or summation for feature fusion to mitigate information loss. Then the relations between the image and question features are captured through computing self-attention on the concatenated features to produce the fused features. CMAN outperforms the commonly used image-question fusion methods in medical VQA as shown in our ablation study.

Third, our model demonstrates superior performance on two large medical VQA benchmarks for both close-end and open-end questions. Especially, our improvement on open-end question answering is overwhelming, with 16% and 41% absolute improvements on VQA-RAD and PathVQA, respectively, verifying the effectiveness of our proposed semi-open framework.

## 2 Method

In this section, we present Q2ATransformer, a semi-open structured model for medical VQA. We first give an overview of our model, and then describe our Visual-Question Encoder in Sec. 3.1 and Answer Querying Decoder in Sec. 3.2.

An overview of our proposed Q2ATransformer model is given in Fig. 2. It follows the majority of medical VQA methods to predict answer classes but exploits candidate answer embeddings for the prediction. Q2ATransformer consists of a Visual-Question Encoder and an Answer-Querying Decoder. The Visual-Question Encoder takes a medical image and a clinic-relevant question as the input and outputs a fused feature with both image and question information. It consists of three parts: vision encoder, question encoder, and fusion network. We use Swin transformer as our vision encoder and BERT as the question encoder. For the fusion network, we propose a Cross-modality Attention Network (CMAN) to integrate image and question features. The Answer Querying Decoder takes the fused image-question feature and learnable candidate answer embeddings as the input and outputs the probability of each candidate answer. Our Answer Querying Decoder consists of two layers of transformer decoders and a classifier to make predictions.

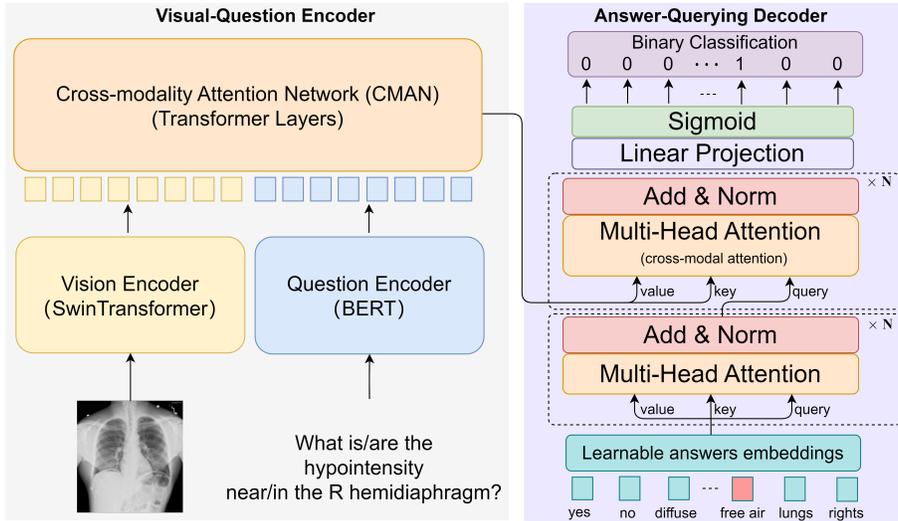
### 2.1 Visual-Question Encoder

The Visual-Question Encoder consists of an image encoder, a question encoder, and a feature fusion module, elaborated as follows.

**Image Encoder.** Our encoder uses the Swin Transformer [12] rather than CNN-based model as our image feature extractor. The advantages of Swin Transformer are three-fold. First, Swin Transformer makes a vision transformer to a hierarchical structure as CNN, which can make the vision transformer more flexible at various scales and has linear computational complexity with the increase of image size. Second, Swin Transformer considers cross-window connection through window shift to obtain long-range dependencies, which introduces more interactions between grids. Therefore, it can provide more regional features and interactions compared with CNN, which is more suitable for the fine-grained nature of medical images. Third, Swin Transformer was pretrained on a large dataset, so it is a very robust feature extractor. Based on these characteristics, we choose Swin Transformer to encode our input image.

Given an input image  $\mathbf{I} \in \mathbb{R}^{H \times W \times C}$ , where  $C$  is the number of channels and  $H$  and  $W$  stand for image height and width, respectively, the image embeddings  $\mathbf{F}_i \in \mathbb{R}^{N \times D_f}$  can be expressed as  $\mathbf{F}_i = \mathbf{W}_i \times \text{SwinTransformer}(\mathbf{I}) + \mathbf{b}_i$ , where  $\mathbf{W}_i$  and  $\mathbf{b}_i$  are learnable parameters to project the output of Swin Transformer into the same dimension  $D_f$  as the question embeddings. They also provide certain flexibility to adapt Swin Transformer to the datasets in our task. Here  $N$  is the number of the extracted image regional features.

**Question Encoder.** For the input question, we use the pre-trained BERT model [3] as the encoder to extract text features. BERT [3] is a successful NLP model. It incorporates context from both directions of a sentence when embedding questions. It has been applied to question answering tasks with the state-



**Fig. 2.** Overview of Q2ATransformer. The input image-question pair is sent to a Visual-Question Encoder to extract and fuse image and question features. The Visual-Question Encoder consists of a Swin-Transformer-based vision encoder, a BERT-based question encoder, and a proposed Cross-modality Attention Network for feature fusion. The fused feature proceeds to the Answer-Querying Decoder, where the input learnable candidate answer embeddings are utilized as the query to compute the attention map and refined according to the attended fused image-question features to predict the presence of the queried answers.

of-the-art results, and is therefore chosen in our task as the question encoder. The question embeddings  $\mathbf{F}_q \in \mathbb{R}^{M \times D_f}$  is obtained by  $\mathbf{F}_q = \text{BERT}(\mathbf{Q}_e)$ , where  $\mathbf{Q}_e$  denotes the input question and  $M$  is the question feature number and  $D_f$  the question feature dimension.

**Feature Fusion Mechanism.** After the image and question features are extracted, respectively, we propose the Cross-modality Attention Network (CMAN) to fuse the information from these two modalities. As medical images are fine-grained and the visual differences of clinical importance are often subtle, we explore a sophisticated way for feature fusion by investigating the interactions between image regional features and question features. In our proposed fusion module CMAN, we first integrate the image features  $\mathbf{F}_i$  and the question features  $\mathbf{F}_q$  by concatenating them together. Compared with the commonly used matrix multiplication or summation for feature fusion, concatenation could mitigate information loss and facilitate the subsequent computation of image-question interaction in our module. After that, the concatenated features are passed to two transformer encoder layers to calculate the relationship between every pair of image question features through the self-attention mechanism of the Transformer. In this way, we could obtain the fused feature carrying the relation

of image question features with minimal information loss. Mathematically, the fused feature  $\mathbf{F}_f$  is obtained as follows.

$$\begin{aligned}
\mathbf{F}_c &= [\mathbf{F}_i; \mathbf{F}_q] \\
\mathbf{Q}_{F_c} &= \mathbf{W}_q \mathbf{F}_c, \quad \mathbf{K}_{F_c} = \mathbf{W}_k \mathbf{F}_c, \quad \mathbf{V}_{F_c} = \mathbf{W}_v \mathbf{F}_c \\
\mathbf{F}_{att} &= \text{Att}(\mathbf{Q}_{F_c}, \mathbf{K}_{F_c}, \mathbf{V}_{F_c}) = \text{softmax}\left(\frac{\mathbf{Q}_{F_c} \mathbf{K}_{F_c}^T}{\sqrt{d_k}}\right) \mathbf{V}_{F_c} \\
\mathbf{F}_f &= \mathbf{W}_f \mathbf{F}_{att} + \mathbf{b}_f
\end{aligned} \tag{1}$$

Here  $\mathbf{W}_q$ ,  $\mathbf{W}_k$ ,  $\mathbf{W}_v$ ,  $\mathbf{W}_f$ , and  $\mathbf{b}_f$  are learnable parameters, and “;” indicates the concatenation operation. The matrices  $\mathbf{Q}_{F_c}$ ,  $\mathbf{K}_{F_c}$ ,  $\mathbf{V}_{F_c}$  are known as the *query*, *key*, and *value* in self-attention calculation, and here they are the linear transformation of the concatenated feature  $\mathbf{F}_c$ .

## 2.2 Answer Querying Decoder

Given an input image question pair, among a set of answers of interest, our Answer Querying Decoder predicts whether each candidate answer matches the corresponding image question pair and uses the candidate with the highest probability as the final answer. For this purpose, we employ a two-layer transformer decoder followed by a linear projector as our classifier, and introduce a set of learnable candidate answer embeddings together with the fused image-question feature  $\mathbf{F}_f$  as the input of the decoder. Assuming there are  $C$  answer classes in total, we need  $C$  candidate answer embeddings with one-to-one correspondence to the  $C$  answer classes. These answer embeddings, collectively represented by a matrix  $\mathbf{A}$ , are randomly initialised and will be updated during training through a self-attention module, a cross-attention module, and a feed-forward network (FFN) in order. Both the self-attention module and the cross-attention module implement the multi-head self-attention ( $MSA(query, key, value)$ ) but with different *key*, *query*, and *value*. The self-attention module computes the relation between different answer embeddings by using  $\mathbf{A}$  to construct all the *key*, *query*, and *value* matrices. The cross-attention module cares about the relation between the answer embeddings  $\mathbf{A}$  and the fused image-question feature  $\mathbf{F}_f$ . It thus uses the answer embedding  $\mathbf{A}$  as the *query* and the fused image-question feature  $\mathbf{F}_f$  as the *key* and *value* to compute the attention and further updates the answer embeddings by combining the attended image-question features. Mathematically, denoting the answer embeddings at the  $l$ -th layer as  $\mathbf{A}_l$ , it will be updated from the output of the previous layer  $\mathbf{A}_{l-1}$  as follows:

$$\begin{aligned}
\mathbf{A}_l &= MSA(\mathbf{A}_{l-1}, \mathbf{A}_{l-1}, \mathbf{A}_{l-1}) \\
\mathbf{A}_l &= MSA(\mathbf{A}_l, \mathbf{F}_f, \mathbf{F}_f) \\
\mathbf{A}_l &= FFN(\mathbf{A}_l),
\end{aligned} \tag{2}$$

where  $l = 1 \dots L$  and  $L$  is the number of Transformer decoder layers. Through this process, the image-question features are injected into the answer embeddings

and used to refine the latter. The refined  $C$  answer embeddings are sent to the final linear projection layer followed by a sigmoid function  $\sigma(\cdot)$  to predict the probabilities of answer classes. That is:

$$\mathbf{p} = \sigma(\mathbf{W}_A \mathbf{A}_L + \mathbf{b}), \quad (3)$$

where  $\mathbf{W}_A$  and  $\mathbf{b}$  are learnable parameters, and  $\mathbf{p}$  is a vector comprising of  $C$  probabilities corresponding to  $C$  answer classes. The answer class with the highest probability is chosen as the predicted answer.

### 2.3 Loss Function

Medical VQA faces a significant class imbalance problem: Yes/No answer classes are much larger than long open answer classes. In order to address the sample imbalance problem more effectively, we choose a simplified asymmetric loss, which is a variant of focal loss while the hyper-parameter  $\gamma$  is set differently for positive and negative classes, as shown in Eqn. 4:

$$\mathcal{L} = \frac{1}{C} \sum_{c=1}^C \begin{cases} (1 - p_c)^{\gamma_+} \log(p_c), & y_c = 1 \\ (p_c)^{\gamma_-} \log(1 - p_c), & y_c = 0 \end{cases} \quad (4)$$

where  $y_c$  is the ground-truth binary label, indicating if the input image-question pair has the answer class  $c$ , while  $p_c$  is the predicted probability for the class  $c$ . The total loss is computed by averaging this loss over all samples in the training data set. We set the hyper-parameters  $\gamma_+ = 1$  and  $\gamma_- = 4$  by default.

## 3 Experiments And Results

### 3.1 Datasets

We conduct our experiments on two medical VQA benchmarks: VQA-RAD [11] and PathVQA [7], which are described as follows.

**VQA-RAD** is the most commonly used radiology dataset seen to date, containing 315 images and 3515 question-answer pairs, each corresponding to at least one question-answer pair. The types of questions include 11 categories: "anomalies", "properties", "color", "number", "morphology", "organ type", "other", and "section". 58% of the questions are close-end questions and the rest are open-end questions. The images are of the body's head, chest, and abdomen. Manual division of the training and test sets is required. For comparability, we divide the data set according to the MMQ method [4].

**PathVQA** is a dataset for exploring pathology VQA. Images with captions were extracted from digital resources (electronic textbooks and online libraries). Open-end questions account for 50.2% of all questions. For the closed-end yes/no questions, the answers are balanced with 8,145 yes and 8,189 no questions. PathVQA consists of 32,799 question-answer pairs, 1,670 pathology images collected from two pathology textbooks, and 3,328 pathology images collected from the PEIR digital library [1]. For comparability, we also divide the data set according to the MMQ method [4].

### 3.2 Comparison with the state-of-the-art methods

We compare our proposed model with 7 state-of-the-art (SOTA) Medical VQA approaches, including StAn [7], BiAn [7], MAML [6], MEVF [13], MMQ [4], PubMedCLIP [5], and MMBERT [9]. The first 6 methods are classification-based approaches. They are chosen because they are among the best performers on the two benchmarks VQA-RAD and PathVQA. The last method MMBERT [9] is chosen as a representative of generation-based approaches, which has the reported performance on VQA-RAD. Except PubMedCLIP [5] and MMBERT [9] whose results are quoted from their original papers, the results of other methods are quoted from MMQ [4]. It is noted that same as our Q2ATransformer, PubMedCLIP [5] and MMBERT [9] employ the same data split as MMQ [4]. Therefore these results are strictly comparable.

As shown in Table 1, on both datasets, our Q2ATransformer consistently outperforms the compared models. Specifically, compared with the second best performer, on VQA-RAD, we achieve an accuracy of 79.19% (16.09% absolute improvement) on Open-ended questions, 81.2% (1.2% absolute improvement) on Close-ended questions, and 80.48% (8.48% absolute improvement) across all questions; on PathVQA, we achieve an accuracy of 54.85% (41.45% absolute improvement) on open-ended questions, 88.85% (4.85% absolute improvement) on Yes/No questions, and 74.61% (25.81% absolute improvement) across all questions. The results could be even better if we increase the dimension of the candidate answer embeddings, as shown in our ablation experiments. From these results, we can see our Q2ATransformer demonstrates overwhelming advantages on open-ended questions, which supports our analysis that by interacting answer information with fused image-question features, our model could better tackle long answer questions. Our model also outperforms the generation-based method MMBERT [9], since we reduce the search space of answers while MMBERT [9] could generate non-existent answers.

**Table 1.** Performance comparison of different methods. † and ‡ indicate the methods are classification-based(closed-type) or generation-based(open-type), respectively.

References Methods	Fusion Methods	PathVQA			VQA-RAD		
		Free-form	Yes/No	Over-all	Open-ended	Close-ended	Over-all
StAn <sup>†</sup> [7]	SAN	1.6	59.4	30.5	24.2	57.2	44.2
BiAn <sup>†</sup> [7]	BAN	2.9	68.2	35.6	28.4	67.9	52.3
MAML <sup>†</sup> [6]	SAN	5.4	75.3	40.5	38.2	69.7	57.1
	BAN	5.9	79.5	42.9	40.1	72.4	59.6
MEVF <sup>†</sup> [13]	SAN	6.0	81.0	43.6	40.7	74.1	60.7
	BAN	8.1	81.4	44.8	43.9	75.1	62.7
MMQ <sup>‡</sup> [4]	SAN	11.2	82.7	47.1	46.3	75.7	64.0
	BAN	13.4	84.0	48.8	53.7	75.8	67.0
PubMedCLIP <sup>†</sup> [5]	-	-	-	-	60.1	80	72.1
MMBERT <sup>†</sup> [9]	-	-	-	-	63.1	77.9	72.0
Ours		<b>54.85</b>	<b>88.85</b>	<b>74.61</b>	<b>79.19</b>	<b>81.2</b>	<b>80.48</b>

**Table 2.** Ablation Studies. BAN, SAN, and CMAN stand for Bilinear Attention Network [10], Stacked Attention Network [16] and ours Cross-modality Attention Network, respectively; Decoder refers to our Answer-Querying Decoder.

#	BAN	SAN	CMAN	Decoder	VQA-RAD			PathVQA		
					open	closed	overall	free-form	yes/no	overall
1	✓				43.62	75.56	64.1	15.03	78.24	51.69
2	✓			✓	54.36	80.07	70.84	44.78	88.29	70.09
3		✓			61.07	77.07	71.33	44.58	86.29	68.88
4		✓		✓	73.83	80.08	77.83	52.88	88.44	73.51
5			✓		69.13	76.32	73.73	47.53	86.73	70.31
6			✓	✓	79.19	81.2	80.48	54.85	88.85	74.61

### 3.3 Ablation study

To investigate the contributions of our proposed feature fusion module CMAN and the decoder for answer querying, we conduct extensive ablation studies to compare different configurations of our model, as presented in Table. 3.2. Here BAN, SAN, and CMAN are three attention networks to fuse image and question features, representing Bilinear Attention Network [10], Stacked Attention Network [16] and ours Cross-modality Attention Network, respectively; Decoder represents our Answer-Querying Decoder. The symbol ✓ indicates the inclusion of the corresponding component. All the experiments in Table. 3.2 are performed based on the same image and question encoders.

**Impact of the CMAN.** The benefit of using CMAN can be well reflected by the improvement from #1 to #5 or from #3 to #5 in Table. 3.2, indicating the effectiveness of our proposed CMAN over BAN and SAN for image and question feature fusion. This is because compared with BAN which multiplies the image and question features or SAN which does a direct matrix summation for fusion, our CMAN directly concatenates the two channels of features together and then calculates attention for fusion. Through this, our CMAN mitigates the information loss due to the multiplication or summation operation during feature fusion in BAN or SAN.

**Contribution of the decoder.** As shown, the inclusion of our answer querying decoder could boost the model performance. To verify the robustness of our decoder, we incorporate it with three different attention modules shown in #2, #4 and #6. By comparing #2 to #1, #4 to #3, or #6 to #5, it can be observed that our answer querying decoder can bring significant performance gain with all three attention mechanisms. Especially, when combining our CMAN and decoder, we can achieve the new SOTA results.

**Impact of answer embedding size.** The experimental results in Fig.3 show that as the dimension of answer embedding increases, the model’s performance improves while the best result is obtained when the embedding size is around 2048. However, increasing the embedding size will also increase the computa-

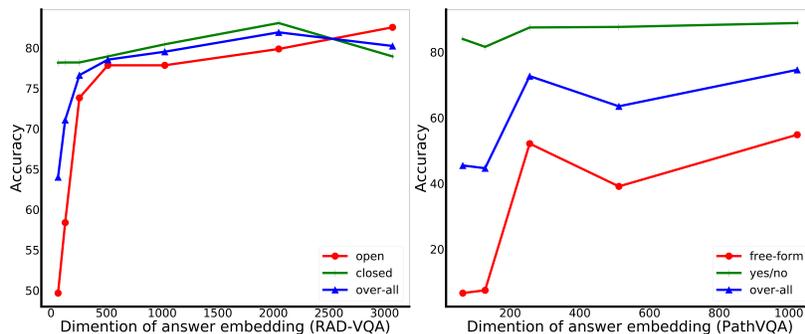


Fig. 3. Ablation study about different dimensions of answer embeddings.

tional cost, while the performance improvement becomes saturated. As a trade-off, our model adopts 1024-dimensional answer embeddings.

### 3.4 Qualitative results

Example results from PathVQA and VQA-RAD datasets are given in Figure 4 and Figure 5, respectively. As can be seen, for these examples where MMQ using BAN for feature fusion fails, our Q2ATransformer w/o decoder has been able to correct most of them using the proposed CMAN fusion module. The performance could be further improved with our Answer Querying Decoder by learning candidate answer embedding through their interactions with the fused image-question features.

### 3.5 Limitation and Discussion

As described in Section 2, we treat each answer as a learnable embedding and use all embeddings as the query to compute the attention map in our decoder. Since we compute the global self-attention, this may increase computation overhead when the number of answer classes is very large. This problem has been encountered in NLP when processing long sequences. Some solutions have been proposed, such as dynamically computing sparse attention, which can significantly reduce computational overhead and will be explored in our future work.

## 4 Conclusion

In this paper, we propose a semi-open framework for medical VQA, which successfully enrolls answer semantic information into the answer class prediction process through our designed mechanism to correlate the answering embeddings with the fused image-question features, which improves the accuracy significantly. It enriches the existing closed-type and open-type medical VQA frameworks and refreshes the SOTA performance on the two benchmarks, especially for the open-end questions.

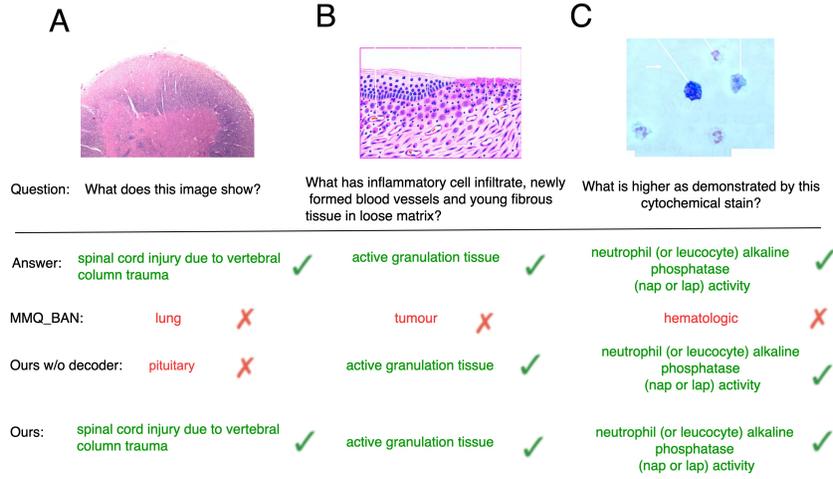


Fig. 4. Example results from PathVQA dataset.

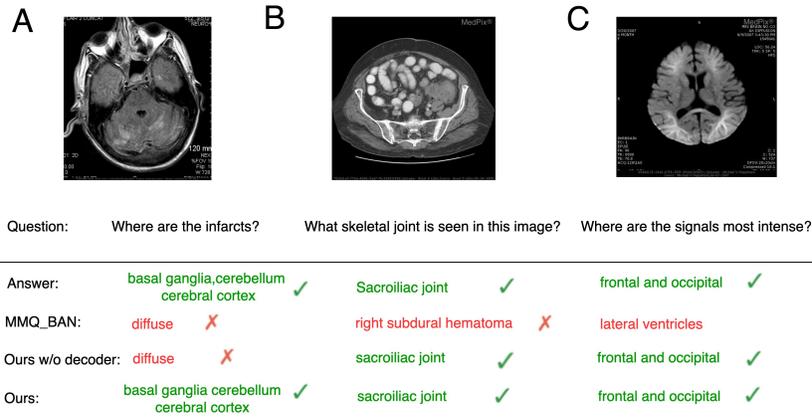


Fig. 5. Example results from VQA-RAD dataset.

## References

1. Peir digital library. <http://peir.path.uab.edu/library/index.php?/category/2>
2. Ambati, R., Dudyala, C.R.: A sequence-to-sequence model approach for imageclef 2018 medical domain visual question answering. In: 2018 15th IEEE India Council International Conference (INDICON). pp. 1–6. IEEE (2018)
3. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)
4. Do, T., Nguyen, B.X., Tjiputra, E., Tran, M., Tran, Q.D., Nguyen, A.: Multiple meta-model quantifying for medical visual question answering. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 64–74. Springer (2021)
5. Eslami, S., de Melo, G., Meinel, C.: Does clip benefit visual question answering in the medical domain as much as it does in the general domain? arXiv preprint arXiv:2112.13906 (2021)
6. Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. In: International conference on machine learning. pp. 1126–1135. PMLR (2017)
7. He, X., Zhang, Y., Mou, L., Xing, E., Xie, P.: Pathvqa: 30000+ questions for medical visual question answering. arXiv preprint arXiv:2003.10286 (2020)
8. Jiang, H., Misra, I., Rohrbach, M., Learned-Miller, E., Chen, X.: In defense of grid features for visual question answering. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10267–10276 (2020)
9. Khare, Y., Bagal, V., Mathew, M., Devi, A., Priyakumar, U.D., Jawahar, C.: Mm-bert: multimodal bert pretraining for improved medical vqa. In: 2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI). pp. 1033–1036. IEEE (2021)
10. Kim, J.H., Jun, J., Zhang, B.T.: Bilinear attention networks. *Advances in neural information processing systems* **31** (2018)
11. Lau, J.J., Gayen, S., Ben Abacha, A., Demner-Fushman, D.: A dataset of clinically generated visual questions and answers about radiology images. *Scientific data* **5**(1), 1–10 (2018)
12. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10012–10022 (2021)
13. Nguyen, B.D., Do, T.T., Nguyen, B.X., Do, T., Tjiputra, E., Tran, Q.D.: Overcoming data limitation in medical visual question answering. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 522–530. Springer (2019)
14. Ren, F., Zhou, Y.: Cgmvcqa: A new classification and generative model for medical visual question answering. *IEEE Access* **8**, 50626–50636 (2020)
15. Wu, C., Liu, J., Wang, X., Li, R.: Differential networks for visual question answering. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 33, pp. 8997–9004 (2019)
16. Yang, Z., He, X., Gao, J., Deng, L., Smola, A.: Stacked attention networks for image question answering. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 21–29 (2016)