

Efficient Bilinear Attention-based Fusion for Medical Visual Question Answering

Zhilin Zhang¹, Jie Wang³, Zhanghao Qin², Ruiqi Zhu³ and Xiaoliang Gong^{3*}

¹Tandon School of Engineering, New York University, New York, USA

²School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore

³College of Electronic and Information Engineering, Tongji University, Shanghai, China
zz10068@nyu.edu, ZHANGHAO001@e.ntu.edu.sg, {2054310, zhurq, gxllshsh}@tongji.edu.cn

Abstract—Medical Visual Question Answering (MedVQA) has attracted growing interest at the intersection of medical image understanding and natural language processing for clinical applications. By interpreting medical images and providing precise answers to relevant clinical inquiries, MedVQA has the potential to support diagnostic decision-making and reduce workload across various fields like radiology. While recent approaches rely heavily on unified large pre-trained Visual-Language Models, research on more efficient fusion mechanisms remains relatively limited in this domain. In this paper, we introduce a fusion model, OMniBAN, that integrates Orthogonality loss, Multi-head attention, and a Bilinear Attention Network to achieve high computational efficiency as well as solid performance. We conduct comprehensive experiments and demonstrate how bilinear attention fusion can approximate the performance of larger fusion models like cross-modal Transformer. Our results show that OMniBAN requires fewer parameters (approximately 2/3 of Transformer-based Co-Attention) and substantially lower FLOPs (approximately 1/4), while achieving comparable overall performance and even slight improvements on closed-ended questions on two key MedVQA benchmarks. This balance between efficiency and accuracy suggests that OMniBAN could be a viable option for real-world medical image question answering, where computational resources are often constrained.

Index Terms—Medical Visual Question Answering, Cross-modal Interaction, Multi-modal Fusion, Bilinear Attention

I. INTRODUCTION

Medical Visual Question Answering (MedVQA) is an emerging field within multi-modal artificial intelligence that adapts the principles of general Visual Question Answering (VQA) to meet the specific demands of the medical domain. The primary goal of MedVQA is to support healthcare professionals by automatically generating accurate answers to clinical questions based on medical images, thereby assisting in clinical decision-making and relieving workload. This task involves the fusion of computer vision and natural language processing techniques to analyze visual data alongside natural language questions to enable contextually relevant and clinically accurate responses.

Compared to general domain Visual Question Answering, MedVQA presents unique and significant challenges. At the image level, medical images such as those in radiology or pathology often exhibit subtle differences within small regions, where even minor pixel variations can represent completely

different diagnostic findings (e.g., tiny lesions). At the text level, clinical inquiries are characterized by specialized terminology and complex language structures. This necessitates language models equipped with domain-specific knowledge to accurately interpret the semantics of medical questions. Given these challenges, the multi-modal fusion module is important, as it must ensure the effective integration of information from these distinct modalities without losing key information.

Despite the recent success of multimodal fusion techniques in enhancing MedVQA performance, there is a significant gap in research focusing on computationally efficient fusion methods. Transformer-based models, especially cross-modal Transformers, have demonstrated strong fusion capabilities and have been widely adopted in this domain. However, these large unified models come with substantial computational demands, making them less suitable for real-time clinical applications or small MedVQA model training, where computational efficiency is essential. Also, a common training strategy for MedVQA involves directly utilizing the embeddings from frozen visual and textual encoders, therefore, the primary computational cost of training concentrates on the fusion network. This motivates the need for exploring alternative fusion techniques that maintain high performance while reducing computational complexity.

In this paper, we propose an efficient fusion framework called OMniBAN, which combines Orthogonality loss, Multi-head attention, and a Bilinear Attention Network. Our approach is designed to deliver comparable performance on MedVQA tasks at a lower training and computational cost. Through extensive experiments, we show that OMniBAN achieves similar results compared with large Transformer-based fusion models on key MedVQA benchmarks, yet it requires fewer computational resources. This efficiency–accuracy trade-off positions OMniBAN as a promising solution for real-world medical image question answering, particularly in radiological Visual Question Answering.

II. RELATED WORK

A. Medical Visual Question Answering

Medical Visual Question Answering (MedVQA) is an emerging research area within multimodal artificial intelligence that applies the general principles of Visual Question Answering (VQA) to the medical domain. This field combines

*Corresponding author.

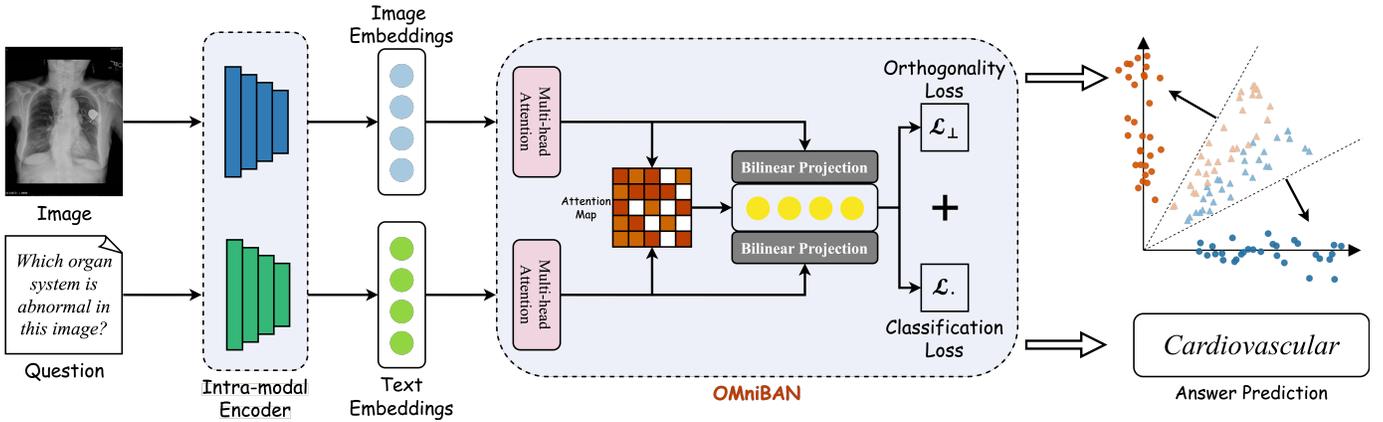


Fig. 1. Overview of our proposed Orthogonal Multi-head Bilinear Attention Network (OMniBAN). The frozen visual and textual backbones encode the input image and question independently, and then the core OMniBAN fusion module fuses these features using multi-head attention for within-modality refinement and bilinear attention to capture cross-modal interactions. Orthogonality loss is adopted to encourage diverse attention patterns among glimpses during training.

medical image understanding and natural language processing techniques to analyze and understand medical images in conjunction with natural language questions, with the goal of generating accurate answers useful for clinical decision-making and diagnostic assistance.

Initial research efforts in MedVQA adopt models that have proven effective in general VQA domain, and adapt them for medical applications. In terms of image feature extraction, researchers often rely on pre-trained models like VGGNet [1] and ResNet [2], which are fine-tuned for the specific task of MedVQA. On the text side, GRU [3] and LSTM [4] are commonly used to extract textual features, while some approaches incorporate additional semantic information derived from medical corpora to enhance the embeddings used for question representation.

To address the challenges specific to MedVQA, such as the scarcity of labeled medical data, various techniques have been proposed. Nguyen et al. [5] introduced the Model-Agnostic Meta-Learning (MAML) framework combined with a Convolutional Denoising Auto-Encoder (CDAE) to improve feature learning. Similarly, Liu et al. [6] utilized contrastive learning to train a pre-trained model (CPRD) that was then applied to MedVQA tasks. These approaches often use transfer learning to leverage external datasets and pre-trained models to enhance the quality of the extracted image and text features.

B. Multi-modal Fusion

In the Visual Question Answering task, multi-modal fusion plays a crucial role by integrating visual and textual features to enable accurate classification. The performance of VQA models largely depends on the effective intra-modal feature extraction and the subsequent inter-modal fusion.

Recent approaches to multi-modal fusion in VQA have introduced various methods to enhance the integration of visual and textual features. A foundational method is SAN [7], which iteratively refines attention on relevant image regions based on the given question. Bilinear pooling has been a major focus in improving fusion by capturing complex interactions

between modalities. Yu et al. [8] proposed MFB to address the high computational cost of bilinear pooling by factorizing the interaction into two low-rank matrices, preserving performance while reducing complexity. MUTAN [9] builds on this by applying Tucker decomposition to further compress the bilinear tensor and get a more compact and efficient multi-modal representation. To make full use of bilinear attention maps, Kim et al. [10] proposed BAN, which can capture dependencies between visual and textual features efficiently.

The advent of Transformer [11] has brought a significant shift in multi-modal fusion strategies. Methods like LXMERT [12], hi-VQA [13], MCAN [14] and METER [15] utilize cross-modality Transformers with separate encoders for vision and language, followed by a cross-modality encoder to fuse the extracted features. In the domain of Medical Visual Question Answering, Liu et al. [16] employs a Transformer-based architecture that directly fuse image and text features to generate joint representation. While Transformers are effective for fusion, they require more computational resources and are more complex compared to bilinear pooling-based methods. This makes balancing performance and computational cost particularly important in Medical Visual Question Answering, where medical data is often limited.

III. METHOD

A. Problem Formulation

Medical Visual Question Answering is regarded as a classification task, and the objective is to identify the most probable answer a from a predefined set of possible answers $A = \{a_1, a_2, a_3, \dots, a_n\}$. This prediction can be expressed as:

$$\hat{a} = \arg \max_{a \in A} P(a | v_i, q_i) \quad (1)$$

where $P(a | v_i, q_i)$ denotes the probability of a given answer a being correct given the image v_i and the question q_i , and \hat{a} is the predicted answer that maximizes this probability.

B. Multi-modal Feature Extraction

Image Encoder. MedVQA requires highly specialized image encoders capable of capturing the intricate details specific to medical images, which often differ significantly from general image data. Medical imaging tasks demand a high level of precision, as even subtle visual cues can hold crucial clinical significance.

In this work, we employ the pre-trained BiomedCLIP Image Encoder [17] rather than other methods as our image encoder. The advantage of BiomedCLIP are two-fold. First, BiomedCLIP builds on the CLIP model [18], which was originally designed to learn image and text representations within a shared feature space through natural language supervision. This design has been proved strong zero-shot performance across various domains. Second, BiomedCLIP further fine-tunes CLIP on the PMC-15M dataset, which consists of diverse medical images and associated text, thereby improving its ability to handle medical visual data.

Given an input radiological image $I_i \in \mathbb{R}^{H \times W \times C}$, BiomedCLIP first produces a hidden representation v_{hid} :

$$v_{hid} = \text{BiomedCLIP}(I_i) \quad (2)$$

This hidden representation is then passed through a projection layer to yield the final 512-dimensional image feature vector v_i representing the initial visual features:

$$v_i = \text{Proj}(v_{hid}) \quad (3)$$

Question Encoder. In this work, We adopt the pre-trained BioBERT model [19] as our question encoder due to its suitability for processing complex biomedical language. BioBERT, based on the BERT architecture, was fine-tuned on a large and diverse biomedical corpus, and demonstrated strong ability in domain-specific biomedical language representation. Therefore, it can generate more accurate representations of medical questions than general-purpose language models.

Compared to traditional VQA models that often use recurrent neural networks, such as LSTM [4] and GRU [3], for text encoding, BioBERT and other BERT-based models provide notable advantages. The Transformer architecture in BERT is particularly effective at capturing long-range dependencies and contextual relationships within text, which can produce richer and more contextually accurate representations of questions. This capability is essential for handling the sophisticated language requirements of MedVQA, where accurate understanding of medical terminology and context is crucial. In our approach, BioBERT encodes each question as a 768-dimensional vector, denoted as q_i .

While BiomedCLIP also offers a text encoder, we chose not to use it for question encoding in this work. Although the BiomedCLIP Text Encoder aligns image and text features within a shared feature space, it lacks the word-level granularity necessary for capturing detailed linguistic information. This level of detail is critical for the fusion methods introduced in Section III-C, which benefit from precise, word-level representations.

C. Orthogonal Multi-head Bilinear Attention Network

Our proposed Orthogonal Multi-head Bilinear Attention Network (OMniBAN) integrates a single-layer multi-head self-attention mechanism with bilinear attention networks to efficiently fuse visual and textual features for Medical Visual Question Answering. This design can help capture complex intra-modal and cross-modal interactions effectively while maintaining computational efficiency. By leveraging orthogonal multi-head attention, OMniBAN enhances feature diversity and maximizes information extraction across modalities.

1) *Intra-modal Feature Attention:* In OMniBAN, we employ a single layer of multi-head self-attention to act as intra-modal attention to refine image and question features independently before cross-modal fusion. Given image features $v_i \in \mathbb{R}^{N_v \times d_v}$ and question features $q_i \in \mathbb{R}^{N_q \times d_q}$, where $N_v = 1$ for image features (since CLIP outputs global image-level features) and N_q denotes the sequence length for question features, we apply multi-head self-attention separately to each modality.

- **Linear Transformations for Queries, Keys, and Values:** For each modality's input x (either v_i or q_i), we generate queries Q , keys K , and values V through linear transformations:

$$Q = xW^Q, \quad K = xW^K, \quad V = xW^V \quad (4)$$

- **Scaled Dot-Product Attention:** We compute attention scores by taking the dot product of Q and K scaled by $\sqrt{d_k}$, and applying a softmax to emphasize relevant information:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (5)$$

- **Multi-Head Attention Output:** The outputs from multiple attention heads are concatenated and linearly transformed to form the final refined features:

$$\tilde{x} = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \quad (6)$$

2) *Cross-modal Bilinear Attention:* After intra-modal refinement, OMniBAN applies a bilinear attention mechanism to fuse the refined visual and textual features. This mechanism captures interactions between modalities by evaluating attention distributions across all pairs of input channels, enabling comprehensive cross-modal feature integration.

To compute the bilinear attention map \mathbf{A} , we use learnable projection matrices $\mathbf{W}_v \in \mathbb{R}^{d_v \times d_m}$ and $\mathbf{W}_q \in \mathbb{R}^{d_q \times d_m}$ for the refined image and question features, respectively, where d_m is the dimension of the shared attention space. The attention map \mathbf{A} is calculated as:

$$\mathbf{A} = \text{softmax}((\mathbf{W}_v \tilde{v}_i) \circ (\mathbf{W}_q \tilde{q}_i)) \quad (7)$$

where \circ denotes the Hadamard (element-wise) product, allowing for efficient alignment between corresponding elements in the image and question features.

The bilinear attention features for each attention head h are then computed by summing over all interactions between image and question features, weighted by the attention map:

$$\mathbf{f}_h = \sum_{j=1}^{N_v} \sum_{k=1}^{N_q} \mathbf{A}_{jk} \left(\tilde{v}_i^j \right)^T \mathbf{W}_{v,h} \mathbf{W}_{q,h} \tilde{q}_i^k \quad (8)$$

In this formulation, each attention glimpse h learns specialized cross-modal relationships to capture diverse interaction patterns between visual and textual features.

3) *Orthogonality Loss*: To ensure diverse information captured by the model, we introduce an **Orthogonality Loss** [20] to encourage each attention glimpse to focus on unique aspects of the input. This loss reduces redundancy across glimpses, which is beneficial considering the subtle and often highly similar pixels present in radiological or pathological medical images. Additionally, it helps address with over-fitting during training, which is common in the field of MedVQA.

In detail, we directly apply Orthogonality Loss (\mathcal{L}_\perp) to the attention distributions obtained from the bilinear attention, as shown in Algorithm 1. For each pair of attention distributions (glimpses), the inner product of their normalized vectors is computed and squared, with these values summed to form the final orthogonality loss.

Algorithm 1 Attention Maps with Orthogonality Loss

Require: Visual features $\mathbf{V} \in \mathbb{R}^{B \times N \times D_v}$, Textual features $\mathbf{Q} \in \mathbb{R}^{B \times T \times D_q}$, Number of glimpses G , Mask M

- 1: **Step 1: Compute Attention Scores**
 - 2: Calculate attention scores: $\mathbf{S} \leftarrow f(\mathbf{V}, \mathbf{Q})$
 - 3: **if** M is applied **then**
 - 4: Apply mask M to \mathbf{S} to ignore invalid regions
 - 5: **end if**
 - 6: **Step 2: Normalize Scores into Distributions**
 - 7: Compute attention distributions: $\mathbf{P} \leftarrow \text{Softmax}(\mathbf{S})$
 - 8: **Step 3: Compute Orthogonality Loss**
 - 9: $\mathcal{L}_\perp \leftarrow 0$
 - 10: **for** each pair of glimpses (g_1, g_2) where $g_1 \neq g_2$ **do**
 - 11: Normalize distributions \mathbf{P}_{g_1} and \mathbf{P}_{g_2}
 - 12: Compute inner product: $\rho \leftarrow \mathbf{P}_{g_1} \cdot \mathbf{P}_{g_2}$
 - 13: Update orthogonality loss: $\mathcal{L}_\perp \leftarrow \mathcal{L}_\perp + \rho^2$
 - 14: **end for**
-

4) *Classifier and Loss Function*: The joint representation output from BAN is then fed into a classifier to predict the most probable answer. A simple feed-forward neural network is used as the classifier, consisting of two fully connected layers with an intermediate activation function.

The main loss function for this task is Binary Cross-Entropy with Logits Loss, which is a common choice for multi-label classification. The total loss during training combines the main classification loss (\mathcal{L}_\bullet) with the Orthogonality Loss (\mathcal{L}_\perp), which encourages both accurate predictions and diversified attention features.

$$\mathcal{L} = \mathcal{L}_\bullet + \alpha \cdot \mathcal{L}_\perp \quad (9)$$

where α is the linear threshold for Orthogonality Loss.

5) *Theoretical Complexity Analysis*: Let $\mathbf{V} \in \mathbb{R}^{N_v \times d_v}$ represent the visual features extracted from N_v image regions ($N_v = 1$ in the context of CLIP-based visual backbones), and let $\mathbf{Q} \in \mathbb{R}^{N_q \times d_q}$ denote the textual features for N_q tokens. We analyze and compare the computational complexity of two fusion paradigms: a vanilla Co-Attention mechanism [14] [21] (which is implemented as a stack of Transformer layers) and our proposed OMniBAN framework.

Both Transformer-based Co-Attention and OMniBAN networks begin by applying intra-modal self-attention to refine the visual and textual features independently. This initial processing step has a computational complexity of:

$$\mathcal{O}_{\text{self-att}} = \mathcal{O}(N_v^2 d_v + N_q^2 d_q) \quad (10)$$

Transformer Co-Attention. A standard Transformer-based co-attention layer processes visual and textual features through cross-modal attention and feed-forward networks (FFN). The dominant operations *per layer* include:

- Cross-modal Attention: Projecting features and performing attention between sequences of length N_v and N_q with dimensions d_v and d_q results in a cost scaling with:

$$\mathcal{O}_{\text{cross-att}} = \mathcal{O}(N_q N_v d_q + N_v d_v d_q + N_q d_q^2) \quad (11)$$

where $\mathcal{O}(N_q N_v d_q)$ denotes the core interaction matrix multiplication cost, assuming that attention dimension is comparable to d_q .

- FFN: For a sequence of length N and dimension d , this costs $\mathcal{O}(N d^2)$. In the co-attention model, this is typically applied to the output sequences of both modalities, and its cost is:

$$\mathcal{O}_{\text{ffn}} = \mathcal{O}(N_q d_q^2 + N_v d_v^2) \quad (12)$$

For a stack of L co-attention layers, the total fusion complexity can be expressed as:

$$\begin{aligned} \mathcal{O}_{\text{co-att}} = & \mathcal{O}(N_v^2 d_v + N_q^2 d_q) \\ & + L \cdot (\mathcal{O}(N_q N_v d_q + N_v d_v d_q + N_q d_q^2 + N_v d_v^2)) \end{aligned} \quad (13)$$

where the dominant terms are:

$$\mathcal{O}(N_v^2 d_v + N_q^2 d_q + L N_q N_v \max(d_v, d_q) + L(N_q d_q^2 + N_v d_v^2)) \quad (14)$$

OMniBAN Fusion. Following the same intra-modal self-attention, OMniBAN employs Factorized Bilinear Attention Networks [10] across γ glimpses. Each glimpse includes:

- Factorized Bilinear Interaction: Computing cross-modal interactions with a cost much lower than $\mathcal{O}(N_v N_q d_v d_q)$ through factorization and projection to a lower intermediate dimension d_m ($d_m \ll \min(d_v, d_q)$), which scales approximately with $\mathcal{O}(N_q N_v d_m)$.
- Post-interaction Projection: A simple projection layer after the bilinear step costs $\mathcal{O}(N_q d_q^2)$.

Therefore, the total fusion complexity is:

$$\mathcal{O}_{\text{OMniBAN}} = \mathcal{O}(N_v^2 d_v + N_q^2 d_q) + \gamma \cdot (\mathcal{O}(N_q N_v d_m) + \mathcal{O}(N_q d_q^2)) \quad (15)$$

Complexity Comparison. Assuming $d_v \approx d_q \approx d$, we can further compare the dominant computational costs of

Transformer-based co-attention and OMniBAN. The initial intra-modal self-attention, which costs $\mathcal{O}((N_v^2 + N_q^2)d)$, is common to both frameworks and less dominant than the costs of the repeated fusion steps for $L, \gamma > 1$. When $L \approx \gamma$, the total complexity difference is primarily determined by the cost per repeated step. OMniBAN’s step is more efficient due to:

- Cross-modal Interaction: Replacing Transformer’s attention $\mathcal{O}(N_q N_v d)$ with factorized bilinear $\mathcal{O}(N_q N_v d_m)$, where $d_m \ll d$.
- Post-interaction Operation: Using a simpler projection $\mathcal{O}(N_q d^2)$ compared to Transformer’s FFN $\mathcal{O}((N_q + N_v)d^2)$. While both exhibit the same asymptotic complexity concerning N_q and d_q , the Transformer’s FFN involves a significantly larger constant factor (due to an expanded intermediate layer), which may lead to higher practical FLOPs.

We also empirically confirm these theoretical efficiency gains in Section IV, where we observe that OMniBAN has fewer parameters and reports much smaller FLOPs compared with co-attention.

IV. EXPERIMENTS

A. Datasets and Metrics

We conduct our experiments on two public medical VQA datasets: VQA-RAD [22] and SLAKE [23]. VQA-RAD contains 3,515 QA pairs derived from 315 radiology images, and is split into 3,064 QA pairs for training and 451 for testing. SLAKE is a Chinese–English bilingual dataset featuring 642 radiology images and a total of 7,033 QA pairs. It includes richer image modalities and covers a broader range of body parts in its questions. For this study, we focus on the English subset of SLAKE (marked as SLAKE-EN in Table II), which consists of 4,919 QA pairs from 450 images in the training set and 1,061 QA pairs from 96 images in the test set. Both VQA-RAD and SLAKE organize their questions into two types: open-ended and closed-ended. Closed-ended questions typically have a limited set of possible responses (most commonly yes/no), while open-ended questions involve more varied answers.

Since MedVQA can be considered as a multi-label classification task, we adopt a binary cross-entropy (BCE) loss function during training. We primarily use accuracy as our evaluation metric, which is computed as the ratio of correctly predicted answers to the total number of questions. To provide a more comprehensive assessment, we report three separate accuracy scores: overall, open-ended, and closed-ended. This breakdown helps compare model performance across different question types.

B. Experimental Setup

We conduct our experiments on a single NVIDIA Tesla V100-SXM2 (16GB) GPU. The learning rate is set to 0.0005, and the batch size is 32. The number of heads in Multi-head Attention and glimpses in BAN are set to 8 and 5, respectively. For the orthogonality loss, we adopt a strategy that linearly increases its weight throughout training (up to

TABLE I
COMPARISON OF ACCURACY (%) ON VQA-RAD [22] TEST SET. THE HIGHEST ACCURACY IN EACH COLUMN IS MARKED IN **BOLD**. RESULTS ARE UNDERLINED TO INDICATE THE IMPROVEMENT BY OMniBAN.

Reference Methods	Fusion Methods	Accuracy		
		Open	Closed	All
MEVF [5]	SAN	40.7	74.1	60.8
MEVF [5]	BAN	43.9	75.1	62.7
MMQ [25]	BAN	53.7	75.8	67.0
CR [24]	BAN	60.0	79.3	71.6
CPRD [6]	BAN	52.5	77.9	67.8
PubMedCLIP(MEVF) [26]	BAN	48.6	78.1	66.5
PubMedCLIP(CR) [26]	BAN	60.1	80.0	72.1
PubMedCLIP(CR) [26]	OMniBAN(Ours)	57.4	<u>80.6</u>	71.4
BiomedCLIP [17]	Transformer	67.6	79.8	75.2
BiomedCLIP(CR) [17]	OMniBAN(Ours)	66.4	<u>80.9</u>	75.1

TABLE II
COMPARISON OF ACCURACY (%) ON SLAKE-EN [23] TEST SET. THE HIGHEST ACCURACY IN EACH COLUMN IS MARKED IN **BOLD**. RESULTS ARE UNDERLINED TO INDICATE THE IMPROVEMENT BY OMniBAN.

Reference Methods	Fusion Methods	Accuracy		
		Open	Closed	All
MEVF [5]	SAN	75.3	78.4	76.5
MEVF [5]	BAN	77.8	79.8	78.6
MMQ [†] [25]	BAN	-	-	-
CR [24]	BAN	78.8	82.0	80.0
CPRD [6]	BAN	79.5	83.4	81.1
PubMedCLIP(MEVF) [26]	BAN	76.5	80.4	78.0
PubMedCLIP(CR) [26]	BAN	78.4	82.5	80.1
PubMedCLIP(CR) [26]	OMniBAN(Ours)	78.1	<u>85.8</u>	<u>81.1</u>
BiomedCLIP [17]	Transformer	82.5	89.7	85.4
BiomedCLIP(CR) [17]	OMniBAN(Ours)	82.0	<u>89.9</u>	85.1

[†]MMQ not reported on the SLAKE dataset.

0.5). We train the models for 40 epochs on both the VQA-RAD dataset and the SLAKE-EN dataset, and save the best-performing model on the validation set as the representative model. Model parameters are optimized using the Adamax optimizer. To mitigate randomness, we train the OMniBAN model ten times with different random seeds and report the average performance along with the standard deviation in Table III.

At the model level, we integrate the CR approach [24], which uses a pre-trained question classifier to determine whether a given question is open-ended or closed-ended. Based on this classification, the question is routed to one of two specialized sub-models. This design directly addresses the distinct linguistic structures and answer formats in open versus closed questions, which can help the overall model capture the subtle semantic differences of question types more effectively. As shown in Table I and II, models that leverage this approach are marked with “CR”.

C. Results and Analysis

Comparisons on VQA-RAD Dataset. Table I shows our results on the VQA-RAD test set, comparing open-ended, closed-ended, and overall accuracy. The baseline models

TABLE III
ABLATION STUDY ON VQA-RAD AND SLAKE-EN TEST SETS (%).

Dataset	Image Encoder	Text Encoder	Multi-head Attention	Orthogonality Loss	Accuracy		
					Open	Closed	All
VQA-RAD	BiomedCLIP [17]	BioBERT [19]	-	-	54.3 ± 3.3	77.3 ± 1.9	68.2 ± 2.0
			✓	-	64.3 ± 1.1	79.4 ± 0.9	73.4 ± 0.5
			✓	✓	66.4 ± 1.0	80.9 ± 1.3	75.1 ± 0.8
SLAKE-EN	BiomedCLIP [17]	BioBERT [19]	-	-	79.0 ± 0.6	84.2 ± 1.4	81.1 ± 1.1
			✓	-	80.6 ± 0.4	87.2 ± 1.2	83.2 ± 0.5
			✓	✓	82.0 ± 0.2	89.9 ± 1.1	85.1 ± 0.6

TABLE IV
COMPARISON OF COMPUTATIONAL EFFICIENCY ON VQA-RAD [22]
TRAINING SET USING PARAMETER SIZES (M) AND FLOPs (M)

Methods	Co-Attention	OMniBAN
Parameters (M)	31.910	21.659
FLOPs (M)	701.276	182.059

MEVF+SAN [5] and MEVF+BAN [5] achieve 60.8% and 62.7% overall accuracy, respectively, with BAN outperforming SAN due to its bilinear attention mechanism. MMQ [25] further raises performance to 67.0%, while CR [24] brings a notable jump to 71.6% by classifying questions into open or closed types before prediction. Leveraging PubMedCLIP with BAN leads to 66.5% accuracy, which increases to 72.1% when combined with CR. Incorporating OMniBAN improves closed-ended performance to 80.6% (PubMedCLIP) and 80.9% (BiomedCLIP), which indicates its strength on questions with restricted answer sets. Although open-ended accuracy dips slightly, overall performance remains competitive, and it demonstrates that OMniBAN can efficiently capture the patterns of closed questions while keeping pace with Transformer-based fusion approaches.

Comparisons on SLAKE-EN Dataset. We observe similar trends on the SLAKE-EN test set in Table II. Again, CR provides a consistent boost over simpler BAN baselines, and OMniBAN excels at closed-ended queries. For instance, PubMedCLIP(CR)+OMniBAN reaches 85.8% on closed-ended questions and 81.1% overall, which improves upon the 82.5% closed-ended accuracy and 80.1% overall accuracy of PubMedCLIP(CR)+BAN. Likewise, BiomedCLIP(CR)+OMniBAN achieves 89.9% closed-ended accuracy, slightly above BiomedCLIP’s 89.7%, but with a marginal decrease in open-ended accuracy.

Efficiency Comparison. In order to demonstrate OMniBAN’s high efficiency, we also compare a typical Transformer-based fusion method, i.e., Co-Attention [14] [21] with OMniBAN, as shown in Table IV. The experiment is conducted on the VQA-RAD training set and evaluated via parameter size and FLOPs. To ensure fairness, the Co-Attention fusion approach includes five layers of cross-modal Transformer encoding (alongside one image–text intra-modal attention layer) to match the five glimpses used by OMniBAN’s bilinear

attention.

The results turn out that OMniBAN requires fewer parameters (21.659M vs. 31.910M) and significantly fewer FLOPs (182.059M vs. 701.276M), which indicates that its bilinear attention component can effectively reduce computational overhead compared with Transformer-level fusion method. This makes OMniBAN a compelling choice for MedVQA scenarios where efficiency is of priority.

D. Ablation Study

We further explore the contribution from each component in the OMniBAN. Table III presents an ablation study on both the VQA-RAD and SLAKE-EN test sets, starting with a baseline model that uses BAN for cross-modal fusion, BiomedCLIP [17] as the image encoder, and BioBERT [19] as the text encoder, without multi-head attention or orthogonality loss. On VQA-RAD, this baseline achieves 54.3% open, 77.3% closed, and 68.2% overall accuracy. Introducing multi-head attention boosts overall accuracy to 73.4%, which reveals the importance of refining intra-modal representations before combining them via BAN. Finally, adding orthogonality loss brings an additional, though smaller, increase to 75.1%. We attribute this limited incremental gain to the inherent constraints of the encoders, which cap the potential benefits of non-overlapping attention glimpses.

On SLAKE-EN, the same pattern emerges. The baseline obtains 79.0% open, 84.2% closed, and 81.1% overall accuracy. Incorporating multi-head attention again yields a notable jump to 83.2% overall, while orthogonality loss provides a further marginal improvement. These findings confirm the value of single-modality attention and orthogonality in enhancing BAN, and suggests that given sufficiently robust features from domain-specific encoders, bilinear fusion can perform competitively with Transformer-based methods for MedVQA.

V. CONCLUSION

In this paper, we introduced the Orthogonal Multi-head Bilinear Attention Network (OMniBAN) as an efficient fusion framework for Medical Visual Question Answering. By combining a single-layer multi-head self-attention mechanism with bilinear attention and employing Orthogonality Loss, OMniBAN balances accuracy and computational efficiency. Experimental results on VQA-RAD and SLAKE-EN show that OMniBAN, when paired with BiomedCLIP, slightly surpasses

the original Transformer-based BiomedCLIP model on closed-type questions and achieves quite similar overall accuracy. This indicates that bilinear attention is capable of capturing the structured patterns often found in such questions, and offers a promising alternative to Transformer-based fusion methods without sacrificing performance.

Beyond efficiency, OMniBAN’s underlying design has broader implications for visual–textual interaction and model adaptability. Future research could delve more deeply into these aspects by examining OMniBAN’s robustness under different data distributions or its ability to integrate external medical knowledge. Refinements to BAN’s internal attention mechanisms also present an avenue for further enhancing cross-modal interactions. Nevertheless, our work has a few limitations. We have not evaluated OMniBAN on broader and more specialized MedVQA datasets (e.g., Surgical VQA), and we have only tested two encoders—PubMedCLIP and BiomedCLIP—which leaves the model’s performance with other potential encoders unexamined. These constraints underscore the need for broader experimentation and more diverse ablations. Overall, our findings show that OMniBAN can serve as an effective and efficient choice for Medical Visual Question Answering, which indicates that bilinear attention deserves continued exploration in medical image understanding.

REFERENCES

- [1] K. Simonyan, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [3] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, “Empirical evaluation of gated recurrent neural networks on sequence modeling,” *arXiv preprint arXiv:1412.3555*, 2014.
- [4] S. Hochreiter, “Long short-term memory,” *Neural Computation MIT-Press*, 1997.
- [5] B. D. Nguyen, T.-T. Do, B. X. Nguyen, T. Do, E. Tjiputra, and Q. D. Tran, “Overcoming data limitation in medical visual question answering,” in *Medical Image Computing and Computer Assisted Intervention–MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part IV 22*. Springer, 2019, pp. 522–530.
- [6] B. Liu, L.-M. Zhan, and X.-M. Wu, “Contrastive pre-training and representation distillation for medical visual question answering based on radiology images,” in *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part II 24*. Springer, 2021, pp. 210–220.
- [7] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola, “Stacked attention networks for image question answering,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 21–29.
- [8] Z. Yu, J. Yu, J. Fan, and D. Tao, “Multi-modal factorized bilinear pooling with co-attention learning for visual question answering,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 1821–1830.
- [9] H. Ben-Younes, R. Cadene, M. Cord, and N. Thome, “Mutan: Multi-modal tucker fusion for visual question answering,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2612–2620.
- [10] J.-H. Kim, J. Jun, and B.-T. Zhang, “Bilinear attention networks,” *Advances in neural information processing systems*, vol. 31, 2018.
- [11] A. Vaswani, “Attention is all you need,” *arXiv preprint arXiv:1706.03762*, 2017.
- [12] H. Tan and M. Bansal, “Lxmert: Learning cross-modality encoder representations from transformers,” *arXiv preprint arXiv:1908.07490*, 2019.
- [13] C. Pellegrini, M. Keicher, E. Özsoy, and N. Navab, “Rad-reconstruct: A novel vqa benchmark and method for structured radiology reporting,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2023, pp. 409–419.
- [14] Z. Yu, J. Yu, Y. Cui, D. Tao, and Q. Tian, “Deep modular co-attention networks for visual question answering,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 6281–6290.
- [15] Z.-Y. Dou, Y. Xu, Z. Gan, J. Wang, S. Wang, L. Wang, C. Zhu, P. Zhang, L. Yuan, N. Peng *et al.*, “An empirical study of training end-to-end vision-and-language transformers,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18 166–18 176.
- [16] Y. Liu, Z. Wang, D. Xu, and L. Zhou, “Q2atransformer: Improving medical vqa via an answer querying decoder,” in *International Conference on Information Processing in Medical Imaging*. Springer, 2023, pp. 445–456.
- [17] S. Zhang, Y. Xu, N. Usuyama, H. Xu, J. Bagga, R. Tinn, S. Preston, R. Rao, M. Wei, N. Valluri *et al.*, “Biomedclip: a multimodal biomedical foundation model pretrained from fifteen million scientific image-text pairs,” *arXiv preprint arXiv:2303.00915*, 2023.
- [18] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [19] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang, “BioBERT: a pre-trained biomedical language representation model for biomedical text mining,” *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020.
- [20] S. Yang, W. Deng, M. Wang, J. Du, and J. Hu, “Orthogonality loss: Learning discriminative representations for face recognition,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 6, pp. 2301–2314, 2020.
- [21] J. Lu, J. Yang, D. Batra, and D. Parikh, “Hierarchical question-image co-attention for visual question answering,” *Advances in neural information processing systems*, vol. 29, 2016.
- [22] J. J. Lau, S. Gayen, A. Ben Abacha, and D. Demner-Fushman, “A dataset of clinically generated visual questions and answers about radiology images,” *Scientific data*, vol. 5, no. 1, pp. 1–10, 2018.
- [23] B. Liu, L.-M. Zhan, L. Xu, L. Ma, Y. Yang, and X.-M. Wu, “Slake: A semantically-labeled knowledge-enhanced dataset for medical visual question answering,” in *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*. IEEE, 2021, pp. 1650–1654.
- [24] L.-M. Zhan, B. Liu, L. Fan, J. Chen, and X.-M. Wu, “Medical visual question answering via conditional reasoning,” in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 2345–2354.
- [25] T. Do, B. X. Nguyen, E. Tjiputra, M. Tran, Q. D. Tran, and A. Nguyen, “Multiple meta-model quantifying for medical visual question answering,” in *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part V 24*. Springer, 2021, pp. 64–74.
- [26] S. Eslami, C. Meinel, and G. De Melo, “Pubmedclip: How much does clip benefit visual question answering in the medical domain?” in *Findings of the Association for Computational Linguistics: EACL 2023*, 2023, pp. 1181–1193.