



INDIAN STATISTICAL INSTITUTE

BT Road, Kolkata – 700108, India

DISSERTATION

**Enhanced Embedding for Multimodal Medical Visual Question
and Answering**

Submitted in partial fulfillment of the requirements
for the award of the degree of

Master of Technology in Computer Science

Submitted by

Akash Suna

Roll No.: CS2405

Under the Supervision of

Prof. Ujjwal Bhattacharya

Professor

Computer Vision and Pattern Recognition Unit

Indian Statistical Institute

Kolkata, India

2026

CERTIFICATE

This is to certify that the dissertation entitled **Enhanced Embedding for Multimodal Medical Visual Question and Answering** submitted by **Akash Suna** (Roll No. CS2405) to the **Indian Statistical Institute, Kolkata**, in partial fulfillment of the requirements for the award of the degree of **Master of Technology in Computer Science**, is a bonafide record of research work carried out by him under my supervision and guidance.

The contents of this dissertation, in full or in parts, have not been submitted to any other Institute or University for the award of any degree or diploma.

I further certify that Akash Suna has fulfilled all the requirements for the award of the degree of Master of Technology in Computer Science.

Prof. Ujjwal Bhattacharya

Professor,

Computer Vision and Pattern Recognition Unit,

Indian Statistical Institute,

Kolkata 700108, INDIA.

ACKNOWLEDGEMENT

I would like to express my gratitude to Dr. Ujjwal Bhattacharya, my advisor at the Computer Vision and Pattern Recognition Unit of the Indian Statistical Institute, Kolkata, for his unwavering support and inspiration throughout my journey. His knowledge and insight have guided me in conducting this research. I am also grateful to all the faculty members of the Indian Statistical Institute for their continuous support and dedication; it is because of them that I am able to add new dimensions to my research.

Finally, I would like to thank my parents and friends for their unwavering support, which left no stone unturned in making this dissertation possible.

Akash Suna

Roll: CS2405

M.Tech CS, 2nd Year

Indian Statistical Institute

Date: 10/06/2026

ABSTRACT

Visual Answering of questions in the field of Medical which is called as (VqA) has grown as a dominant area of research that fuse processing of natural language and vision of computer often known as CV or NLP to assist in medical decision-making. However, effective multimodal fusion between medical images and clinical questions remains a significant challenge. This thesis examines the application of the Perceiver IO architecture as an efficient multimodal aggregator for medical VQA. The work has been carried out in multiple directions. First, a classification-based framework is developed by combining Vision Transformer (ViT) and ClinicalBERT alongside a Perceiver IO aggregator to perform multimodal fusion for generating answers to medical questions. In the second approach, the florENCE TWO vision-language model is finely tuned with parameter-efficient techniques to enable generative medical VQA. Finally, a hybrid architecture is introduced, where Perceiver IO is employed as a fusion module to integrate visual and textual representations, which are then used to condition the Florence-2 model for answer generation. In this thesis, the VQA-RAD dataset and ImageCLEF Med VQA 2019 dataset are used for the experiments. Through these explorations, the thesis examines both the classification and generative paradigms in area of medical Visual Question and answering and analyzes the importance of Perceiver IO for multimodal understanding. Several challenges encountered in hybrid modeling are discussed, along with right direction in terms of research in future, including the growing development of effective fusion strategies and the extension of the proposed approaches to larger medical VQA datasets. Overall, this thesis contributes to the understanding of efficient multimodal fusion techniques and provides insights for building both classification and more general systems in medical question along with answers on visual data

Contents

1. Introduction	4
1.1 Objectives	5
1.2 Problem Definition	7
2. Existing Techniques and Related Work	8
2.0.1 Traditional Fusion Techniques	8
2.0.2 Deep Learning and Transformer-based Approaches	9
2.0.3 Perceiver IO Architecture	9
2.0.4 Large Vision Language Models	9
2.0.5 Hybrid and Efficient Fusion Approaches	10
2.0.6 Research Gap	10
2.1 Datasets	10
2.1.1 VQA-RAD Dataset	10
2.1.2 Description	11
2.1.3 VQA-Med 2019 Dataset (ImageCLEF)	12
2.1.4 Categories of Questions	13
2.2 Existing Techniques in Question Answering on Medical Data	14
2.2.1 Bilinear techniques for Pooling	14
2.2.2 Transformer Based Architectures	15
2.2.3 Vision-Language Pretraining Models	15
2.2.4 Perceiver-based Fusion Approaches	15
2.3 Evaluation Metrics	16
2.3.1 Accuracy	17
2.3.2 BLEU (Bilingual Evaluation Understudy)	17
2.3.3 METEOR (Metric for Evaluation of Translation with Explicit ORdering)	17
2.3.4 ROUGE-L (Recall-Oriented understudy for gisting evaluation - Longest Common Subsequence)	18
2.3.5 Summary of Evaluation Metrics	18
3. Proposed Method for Medical VQA	20
3.1 The Perceiver IO based Classification Model	20
3.1.1 Vision Transformer (ViT) for Image Feature Extraction	20
3.1.2 Key Features of the Vision Transformer Architecture	22

3.1.3	ClinicalBERT for Clinical Question Understanding	22
3.2	Perceiver IO Architecture	24
3.3	Overall Architecture, Training Procedure, and Results	26
3.4	Fine-tuned Florence-2 Generative Model	27
3.5	Hybrid Architecture: Perceiver IO + Florence-2	29
3.6	Comparative Analysis	31
3.6.1	Failure Analysis and Mitigation Strategies	34
4.	Conclusion	39
5.	Future Work	41
6.	Bibliography	42

Chapter 1

Introduction

Medical imaging has now become an indispensable tool in today's healthcare, helping physicians diagnose and monitor a broad range of diseases with high accuracy. With the massive growth in medical imaging data, there is an increasing demand for intelligent systems that help radiologists incorporate complex visual information. Question Answering which is visual in medical domain often called as Vqa which is growing as a promising area to research upon that fuses language such as natural language and vision of computer to enable machines to show result clinical medical image questions. While VQA is generally considered in terms of general domain-related issues such as scarcity of labeled training data, class overlap, and the requirement for accurate and explainable responses, medical VQA faces unique challenges. In particular, the fusion of visual and textual modalities remains one of the most challenging issues in this area. Conventional approaches to the fusion of both modalities often fall short. While there has been considerable progress made towards multi-modal learning, conventional fusion strategies such as concatenation and attention have limitations in capturing the intricate relations that exist between the visual and textual features in medical VQA. These methods are either computationally intensive or lack the capability to learn long-range dependency features from high-dimensional data.

1.1 Objectives

The key goal of this thesis is to analyze and devise effective multimodal fusion strategies that can be employed to solve the task of question along with answers on Visual data often called as (VqA). Medical VQA seeks to create intelligent algorithms that are able to answer clinical questions posed on the basis of images, which will help radiologists and physicians make informed decisions. Yet, one of the basic problems of Medical VQA relates to the issue of efficiently combining visual information available in the medical images and textual information contained in the clinical questions. To tackle such difficulties, this work investigates how the Perceiver IO system can be used due to its effective processing of high-dimensional multimodal data based on a latent bottleneck structure. The research is organized according to the following specific goals: In the first place, a medical VQA system relying on classification is proposed through the fusion of ViT and ClinicalBERT with the use of Perceiver IO as the multimodal aggregator. Such a method is centered on enhancing learning of joint representations for answering clinical questions correctly. In the second place, the Florence two vision-language model which is finely tuned based on techniques which are efficient on parameter in order to implement generative medical Visual Question and Answering. Finally, a hybrid architecture is suggested and studied, in which Perceiver IO is used as the fusion mechanism. The visual and textual information is combined using Perceiver IO and then fed into the Florence-2 model to assist in answering questions posed to the system. This hybrid model aims to make the best use of both structured multimodal fusion and the generative ability of vision-language models. In evaluating the suggested models, VQA-RAD will be used as the primary test set. However, their usability with other datasets such as VQA-Med 2019 will also be explored.

1.2 Problem Definition

Medicine-based Visual Question Answering (VQA) has a number of distinctive issues different from those of general VQA models. One of the key difficulties in medical VQA is the appropriate fusion of visual and textual input data. Since medical images may be quite complex, and at the same time medical questions require an understanding of specific terminologies, existing fusion algorithms like simple concatenation, element-wise product, or attention mechanisms may be insufficient to detect important semantic connections between these input components. A further important problem is the low performance on questions which are opened and questions which are closed. While closed ones Yes or No"questions have reasonable precision rates, open-ended descriptions have very low accuracies due to the variety of possible responses, class imbalance, and the need for linguistically and contextually relevant answers.



Figure 1.1: Example of a chest X-ray image from which clinical questions are asked in Medical VQA

For instance, with an input like the chest radiograph in Figure 1.1, one might ask the computer system such as

- Is the appearance of the lung normal?: Primarily which is Yes/No type questions.
- Is there any evidence of pleural effusion?: Questions that requires descriptive or detailed answers.

Chapter 2

Existing Techniques and Related Work

Answering into the Questions of Visual Data which is called as "VQA" is a crucial study area which act as hurdle between language of natural processing and vision of Computer which is called as NLP AND CV. The main aim of visual medical question and answering is to create algorithms capable of answering clinical questions with the help of images like x-RayS, CT scANs, and images of MRI. Many techniques have been developed to tackle challenges such as multimodal fusion, domain knowledge incorporation, and answer generation in medical VQA.

2.0.1 Traditional Fusion Techniques

Earliest attempts to develop VQA systems used basic fusion techniques for merging image and text representations. Methods involving concatenation, addition/multiplication operations, or bilinear pooling were commonly adopted. Although these methods are simple to deploy, they lack the capability of identifying more advanced relationship ON text and visual data. To help increase the effectiveness of such fusions, attention mechanisms were introduced which permit the network to attend only to those parts of an image that correlate with the input question. One of the first techniques to do so was stacked attention networks (SANs).

2.0.2 Deep Learning and Transformer-based Approaches

As a result of deep learning’s success, CNNs, which include VGGNet, ResNet, and DenseNet, have emerged as popular models to extract features from images, whereas RNNs and later transformers came into play to help encode the questions. The introduction of transformers has appeared to be groundbreaking for the field. Models such as the Vision Transformer (ViT) and BERT helped enhance representation learning for images and questions. In the case of clinical applications, BioBERT and ClinicalBERT models were introduced that could help represent clinical language. Several studies combined ViT and different variants of BERT in medical VQA tasks.

2.0.3 Perceiver IO Architecture

The IO variant of the Perceiver model is known for its efficiency compared to conventional transformers for handling large multi-modal inputs. Contrary to conventional transformers, which do not scale well with respect to input size, Perceiver IO employs a latent bottleneck for cross-attention mechanisms that enables it to efficiently process large inputs without incurring higher computational overheads. The IO model is particularly relevant for medical visual question answering tasks due to the high-dimensionality of both clinical images and text.

2.0.4 Large Vision Language Models

In more recent years, pre-trained vision-language models like Florence-2, LLaVA, and GPT-4V have shown impressive performance in generating answers to questions based on images. Such models are trained using large image and text datasets and can provide answers without requiring additional processing. In medicine, researchers have attempted to adapt such models to the medical setting by fine-tuning them on medical datasets. The Florence-2 model is particularly interesting because of its unique design that can perform different kinds of vision-language tasks, among which is VQA.

2.0.5 Hybrid and Efficient Fusion Approaches

Many studies have explored hybrid networks incorporating various model architectures. Some researches have utilized fusion blocks that can leverage features from vision and language encoders pre-trained outside their networks. Other studies have examined parameter-efficient adaptation approaches including LoRA that can adapt large models to target domains without full training. However, despite recent progress, performing efficient multimodal fusion for generating answers in medical VQA still poses a challenge. Existing methods tend to adopt basic fusion approaches or demand considerable computational resources for multimodal fusion in VQA tasks.

2.0.6 Research Gap

While there has been considerable success in the development of techniques related to VQA in medical applications, there are still shortcomings with regard to their efficiency, scalability, and versatility in handling multiple tasks including classification and generation in the same system. The need therefore arises to develop systems which are able to efficiently fuse different modalities of medical information while at the same time maintaining low computational requirements. This thesis attempts to partially fill this gap through the use of Perceiver IO fusion in medical VQA.

2.1 Datasets

In this thesis, two publicly available benchmark datasets are used to evaluate the proposed models: VQA-RAD and ImageCLEF Med VQA 2019 (VQA-Med 2019).

2.1.1 VQA-RAD Dataset

The question along with answer about visuals regarding to Radiology often called as VqA Rad dataset can be considered one of the earliest available datasets for Medical Visual Question Answering. This dataset was originated

by Lau et al. in the year 2018 with the aim of overcoming the absence of clinically relevant question along with answers on visual data in the area of medical The VQA-RAD differs from other general-purpose VQA datasets (like VQA v1/v2, GQA) that mostly work on natural scenes and everyday questions. As compared to those datasets, this dataset consists of radiology images only and provides medically relevant questions asked by medical professionals.

Dataset Statistics

The VQA-RAD dataset contains of : Three hundred and fifteen images of radiology and three thousand five hundred and fifteen question along with their answers. The questions about the data are fragmented into following:

- Closed questions: Primarily which is Yes/No type questions.
- Open questions: Questions that requires descriptive or detailed answers.

This dataset is considered to be challenging because of its relatively small size and high diversity in imaging modalities.

2.1.2 Description

The images of the VQA-RAD dataset are sampled from MedPix, an open-access online database of peer-reviewed radiology teaching cases. The three hundred and fifteen images are evenly distributed over three anatomical regions, namely the head (axial CT and MRI scans), the chest (frontal X-rays) and the abdomen (axial CT and MRI scans), so that the dataset covers a wide variety of imaging modalities, organs and pathological conditions. The questions were naturally asked by clinicians while looking at the images and were subsequently validated, and each question is paired with a reference answer. The questions belong to eleven clinically meaningful categories, which include modality, plane, organ system, abnormality, presence of an object or condition, positional reasoning, color, size, attribute, counting and other. The answers are either short free-form phrases in the case of open-ended questions or “yes/no” type responses in the case of closed-ended questions. An official split of the question–answer pairs into a training set and a test set is provided, which enables fair comparison

among different models. Figure 2.1 shows a representative sample from the dataset consisting of a radiology image together with its question and answer pairs, while Figure 2.2 illustrates the wide variety of imaging modalities and anatomical regions present in the dataset.

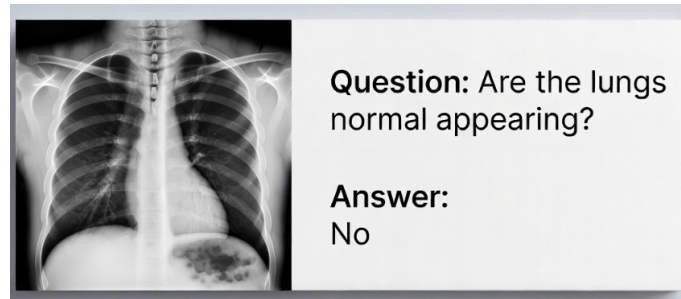
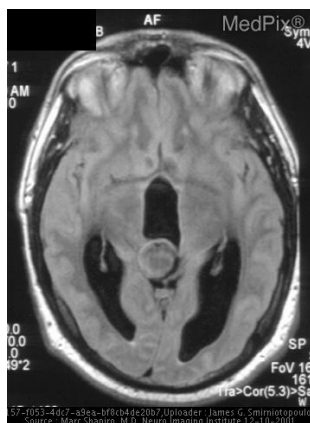


Figure 2.1: Example of what dataset contains an image and the question and answers pair in field of Medical VQA



Q: Which imaging modality was used to acquire this image?

A: CT scan

Q: Is this an axial plane image of the head?

A: Yes



Q: Which organ system is shown in this image?

A: Abdomen

Q: Is this image acquired using MRI?

A: Yes

Figure 2.2: Sample images from the VQA-RAD dataset together with their question and answer pairs: a head CT scan (left) and an abdominal MRI scan (right)

2.1.3 VQA-Med 2019 Dataset (ImageCLEF)

The dataset which is medical data set visual question answer called VQa Med 2019 was proposed as a new benchmark for the question along with answers in the are of medical track of ImageCLEF 2019 competition. The creation of this data set aims to help develop question along with answers in the field

of medical due to its bigger size compared to other data sets like VQA-RAD. Contrary to VQA-RAD that has few images and open-ended questions, the VQA-Med 2019 data set provides a substantially large number of radiology images as well as questions classified according to clinically relevant categories.

Dataset Statistics

This visual question and answers on medical 2019 dataset contains four thousand two hundred radiology images and over 15 thousands question along with answers. These pictures are obtained from diverse imaging methods like XRays, CTScans, and scans of Mri. There exists an official partition between the training and testing set of the database that facilitates comparative analysis of the participating systems within the ImageCLEF challenge.

2.1.4 Categories of Questions

One of the important features of the Visual Question and answers on 2019 data is to do the classification of questions into four broad categories. The inclusion of such categories helps to conduct a detailed analysis of the performance of models in various clinical situations. These four categories are given below which are:

- **Modality:** Questions which is regarding to the imaging modality used for image acquisition (for example, "Which imaging modality has been used?").
- **Plane:** This category includes questions regarding the plane of the image (e.g., "What plane of the image is shown?").
- **Organ:** Questions about recognizing the organ or body part that appears to be in the image.
- **Abnormality:** This category includes questions which is related to abnormality or pathology present inside the image.

This categorization helps in assessing the performance of the models for specific clinical situations.

2.2 Existing Techniques in Question Answering on Medical Data

In the area of giving answers related to the question on visual data called as (VQA) has developed quite extensively over the last few years due to developments in computer vision, NLP, and multimodal learning domains. There have been several approaches formulated in order to merge visual and textual data in order to answer queries related to clinical images. This subsection elaborates on the main types of approaches available in the literature.

1. **Early Fusion Techniques** In the early stages of medical VQA, approaches involved simple fusion mechanisms for the combination of image and text features. Feature concatenation, feature addition, and Hadamard product were some of the most common techniques. In this strategy, the visual features which were obtained with the help of Convolutional Neural Networks like VGGNet or ResNet while the textual features were extracted with the help of word embeddings through Recurrent Neural Networks (RNN) or LSTMs like Word2Vec and GloVe. These mechanisms are very efficient computationally but fall short when capturing intricate relationships between the two types of input features.
2. **Attention based Methods** In order to overcome the limitations of fusion alone, attention was adopted in the medical VQA field. Using attention, it is possible to direct the attention of the model related to the desired parts for an image according to the input query. Some notable approaches include Stacked Attention Networks (SAN) and Co-Attention approach in which the model learns from both visual and linguistic features using an iterative approach. With the help of attention mechanisms, it is now possible to dynamically align the query with the important image regions.

2.2.1 Bilinear techniques for Pooling

The bilinear pooling techniques which were introduced to create the models which can represent more effective interactions among visual and textual fea-

tures. Some popular methods which involve Bilinearity of Pooling which is multimodal called as "McB", Multimodal Pooling of Low rank which is bilinear called as "MLB", and MUTAN, which make use of the final product of the operation. However pooling which are bilinear their techniques is giving better results than other methods for both the generic and medical VQA problem, they are having high computational complexity and memory consumption, which makes them unsuitable for real-world application scenarios.

2.2.2 Transformer Based Architectures

However, the rise of the Transformer model was considered to be major achievement in the development of multimodal learning. Attention-based architectures is useful for more efficient representation of long dependencies between images and texts. In the medical field, the combination of Vision Transformers for representing visual informations along with BERT-like models (ClinicalBERT, BioBERT) for handling questions became common practice. This technique proved to be superior to traditional CNN-RNN pipelines due to its ability to represent more fine-grained concepts.

2.2.3 Vision-Language Pretraining Models

More recently, however, language models for visuals which is often large known as "VLMs" that are already trained with massive scale picture along with text pairs have been used in VQA for the field of medicine. Models like CLIP, VisualBERT, and Florence-2 have performed well in zero and fewshot settings. In the field of medicine, researchers have fine-tuned these VLMs with datasets from the domain of radiology for performing VQA tasks. One model named Florence-2, specifically, has been in the spotlight due to its multi-task capability in vision language tasks.

2.2.4 Perceiver-based Fusion Approaches

In order to overcome the efficiency problem that arises with attention-based and bilinear approaches, the Perceiver IO was proposed as an efficient fusion

architecture for multimodal learning. This model adopts a latent bottleneck which enables the model to iteratively work on high-dimensional inputs, and thus, it can scale up to very large inputs including medical images with high resolutions and lengthy clinical queries. The Perceiver IO has been proved successful in general multimodal tasks, and currently, it is receiving growing attention in medical applications.

Summary

Figure 2.3 presents a summary of the existing techniques in Medical VQA discussed in this chapter.

Technique	Key Idea	Strengths	Limitations	Examples
Early Fusion	Simple concatenation / multiplication	Simple and fast	Weak multimodal interaction	CNN + LSTM + Concatenation
Attention Mechanisms	Dynamic focus on relevant regions	Better alignment	Limited higher-order interactions	SAN, Co-Attention
Bilinear Pooling	Outer product for feature interaction	Rich multimodal representation	High computational cost	MCB, MLB, MUTAN
Transformer-based	Self-attention for sequences	Strong contextual understanding	Requires large data	ViT + BERT, VisualBERT
Vision-Language Models	Pretraining on large image-text data	Strong generalization	High resource requirement	CLIP, Florence-2
Perceiver IO	Latent bottleneck for efficient fusion	Scalable and efficient	Requires careful design	Perceiver IO in Medical VQA

Figure 2.3: Summary of existing techniques in Medical VQA

2.3 Evaluation Metrics

For thoroughly examines the performance of the two model, metrics for classification tasks and metrics for generation tasks were used. The choice of metrics was influenced by the questions type asked in the Vqa Rad data, as it contains both closed and open-ended questions.

2.3.1 Accuracy

Closed-ended questions which were generally found to be yes/no questions were evaluated using the metric called Accuracy. This metric helps to compute the accuracy for predicting the answers to the total number of closed-ended questions. In mathematical terms, Accuracy can be written as:

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Closed-ended Questions}} \times 100 \quad (2.1)$$

For closed ended question this metric is turned out to be very useful because of their limited answer space, which exact matching helps to provide a clear and interpretable measure of performance.

2.3.2 BLEU (Bilingual Evaluation Understudy)

The **BLEU** markings help us to analyse the quality of generated answers by measuring the overlapping for n-Gram among the predicted answer and the ground truth answer. It computes updated nGrams accuracy alongside a brevity penalisation overly short predictions. The scores are calculated as:

$$\text{BLEU}_n = \text{bp} \times \exp \left(\sum_{r=1}^n w_r \log p_r \right) \quad (2.2)$$

where p_r denotes the modified accuracy for r -grams, w_r represents the weight for each r -gram (typically as $1/n$), and "bp" is the brevity penalisation.

2.3.3 METEOR (Metric for Evaluation of Translation with Explicit ORdering)

METEOR is an evaluation metric which helps to evaluate the output text by taking care of both precision and recall at the unigram level while also considering synonymy and stemming. It first computes the unigram precision P and the unigram recall R between the generated answer and the reference answer, which are then combined into a weighted harmonic mean:

$$F_{\text{mean}} = \frac{10 \cdot P \cdot R}{R + 9 \cdot P} \quad (2.3)$$

To account for the ordering of words, a fragmentation penalty is computed based on the number of contiguous matched chunks:

$$\text{Penalty} = 0.5 \times \left(\frac{\text{number of chunks}}{\text{number of matched unigrams}} \right)^3 \quad (2.4)$$

The final METEOR score is then given by:

$$\text{METEOR} = F_{\text{mean}} \times (1 - \text{Penalty}) \quad (2.5)$$

2.3.4 ROUGE-L (Recall-Oriented understudy for gisting evaluation - Longest Common Subsequence)

ROUGE-L helps to evaluate the similarities in the generated results(answer) and the reference answer which is based on the Long common subSequence "LcS". Given a reference answer X of length m and a generated answer Y of length n , the recall and precision based on the longest common subsequence are computed as:

$$R_{\text{lcs}} = \frac{\text{LCS}(X, Y)}{m}, \quad P_{\text{lcs}} = \frac{\text{LCS}(X, Y)}{n} \quad (2.6)$$

The final ROUGE-L score is the F-measure computed from these two quantities:

$$\text{ROUGE-L} = F_{\text{lcs}} = \frac{(1 + \beta^2) R_{\text{lcs}} P_{\text{lcs}}}{R_{\text{lcs}} + \beta^2 P_{\text{lcs}}} \quad (2.7)$$

where $\text{LCS}(X, Y)$ denotes the length of the longest common subsequence between X and Y , and β controls the relative importance assigned to recall over precision.

2.3.5 Summary of Evaluation Metrics

The metrics which are used for evaluation in this thesis are given in Table 2.1.

Table 2.1: Metrics for evaluation

Metric	Question Type	Purpose	Formula
Accuracy	Closed-ended	Measures exact correctness	Eq. 2.1
BLEU	Open-ended	Evaluates n-gram overlap	Eq. 2.2
METEOR	Open-ended	Captures semantic similarity	Eq. 2.5
ROUGE-L	Open-ended	Measures structural similarity via LCS	Eq. 2.7

Chapter 3

Proposed Method for Medical VQA

The question along with answer of visuals in field of medical domain called as "vQa" problem involves the use of the model for information extraction from the radiology image and questions asked on the clinics and then generate answers accordingly. The main difference between the general question answers on visual problem and the medical question answerins on visual problem includes the need for accuracy and domain-specific knowledge as well as handling yes/no and descriptive questions. Three tracks are proposed to overcome the difficulties outlined above.

3.1 The Perceiver IO based Classification Model

In this regard, our proposed architecture consists of 3 components: Visual Feature Extractor (Vision Transformer), Natural Language Processor (ClinicalBERT), and Multimodal Fusion Module (Perceiver IO).

3.1.1 Vision Transformer (ViT) for Image Feature Extraction

ViT acts as the image encoder in Track 1. The Convolutional Neural Network (CNN) has been the dominant architecture in computer vision applications for decades now. CNNs, however, have local receptive fields and reliability of range which is long are hard to find between patches in img. In medical image processing, such dependencies between distant patches might play a critical role in determining the existence of diseases that may have complex spatial relations. On the other hand, ViTs address this problem by dividing the input

of the images in terms of patches and processing each patch through a self-attention mechanism like transformers do with text inputs. Given an image input $x \in \mathbb{R}^{H \times w \times C}$, changing a shape in sequential patches of two dimension which are flattened $X_p \in \mathbb{R}^{n \times (P^2 \cdot C)}$, where (p, P) size of the patch and $n = hW/p^2$ is the count of patches involved. Each and every ptch is project linearly into a D -dimensional embedding space with a learnable matrix E :

$$x_p^i = \text{Flatten}(x^i) \cdot E, \quad i = 1, 2, \dots, N$$

A classification which is learnable token x_{cls} is prepended sequence, and learnable embeddings of position E_{pos} are added to retain information regarding distance lost during flattening. The overall architecture of the adapted Vision Transformer is illustrated in Figure 3.1.

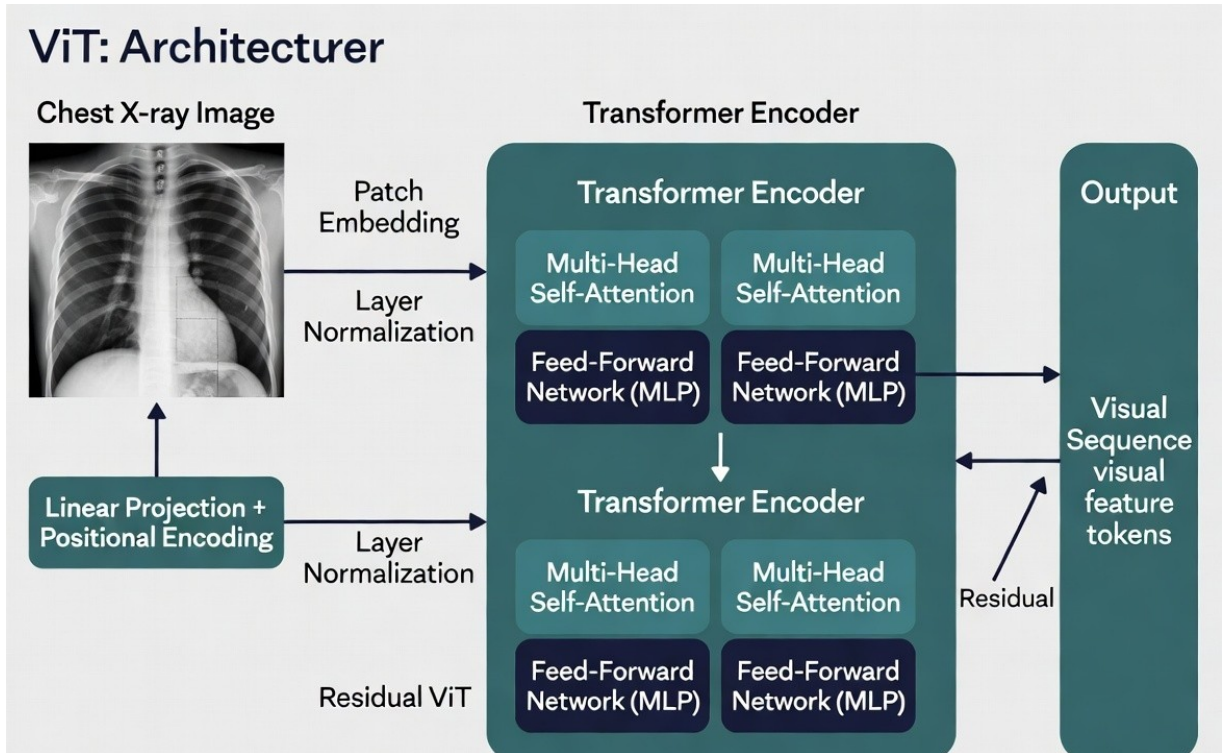


Figure 3.1: Vision Transformer (ViT) architecture adapted for medical imaging: the input radiology image is split into fixed-size patches which are flattened and linearly projected into patch embeddings; positional embeddings and a learnable classification token are then added, and the resulting sequence is processed by a stack of transformer encoder layers (multi-head self-attention, MLP, layer normalization and residual connections) to produce the output visual tokens (adapted from [3])

3.1.2 Key Features of the Vision Transformer Architecture

The introduced architecture related to VISION transFORMER adapted for the medical domain has the following key characteristics:

- **Medical Context:** Unlike standard ViT models trained on natural images, this architecture is specifically designed for medical imaging. It accepts radiology images (such as chest X-rays from the VQA-RAD dataset) as input.
- **Patch Embedding:** The image of medical which is going to be our input is fragmented into a non OverLAP patches sequences. PaTches of This are projected which is linearly into embeddings, enabling the transformer to process the tokens of sequence images.
- **Transformer Encoder:** The core of the model consists of multiple transformer encoder layers. Self-Attention which is of Multi-hEad that each layer is having, network which is feedForwARD (MLP), Layer Normalization, and Residual Connections.
- **Output Representation:** visual tokens which is our final output sequence tokens containing rich semantic information from the medical image.
- **Academic Suitability:** The architecture follows a clean and modular design, making it suitable for thesis and research documentation.

3.1.3 ClinicalBERT for Clinical Question Understanding

To process clinical questions, we use ClinicalBERT, a domain-specific language model obtained via additional fine-tuning of BERT on massive amounts of clinical texts in the MIMIC-III corpus. Generic language models such as BERT are often unable to capture the specifics of medical terminology, abbreviations, specific vocabulary, and clinical reasoning patterns. Like its predecessor, ClinicalBERT is based on the Transformer Encoder. Given a clinical question

$q = \{[\text{CLS}], w_1, w_2, \dots, w_L, [\text{SEP}]\}$, ClinicalBERT produces contextualized token representations:

$$H_q = \{h_{[\text{CLS}]}, h_1, h_2, \dots, h_L, h_{[\text{SEP}]}\} \in \mathbb{R}^{(L+2) \times d}$$

where $d = 768$ for the base model. The embedding which is related to the special [cLs] token ($H_{[\text{cLs}]}$) is used as the combined representation of the overall questions. ClinicalBERT is pre-trained using two objectives: modelling of language which is masked "MIM" and the Prediction of next Sentence "nSP". In the context of Medical VQA, ClinicalBERT provides significant advantages. It can better understand medical questions containing specialized terminology (e.g., "infarcted", "consolidation", "effusion") and capture the clinical intent behind the question. This domain adaptation leads to more accurate multimodal alignment between images and questions. For integration into our proposed architecture, we extract the contextualized token embeddings H_q (excluding the [SEP] token) and pass them to the **Perceiver IO Aggregator**. The overall architecture of ClinicalBERT is illustrated in Figure 3.2.

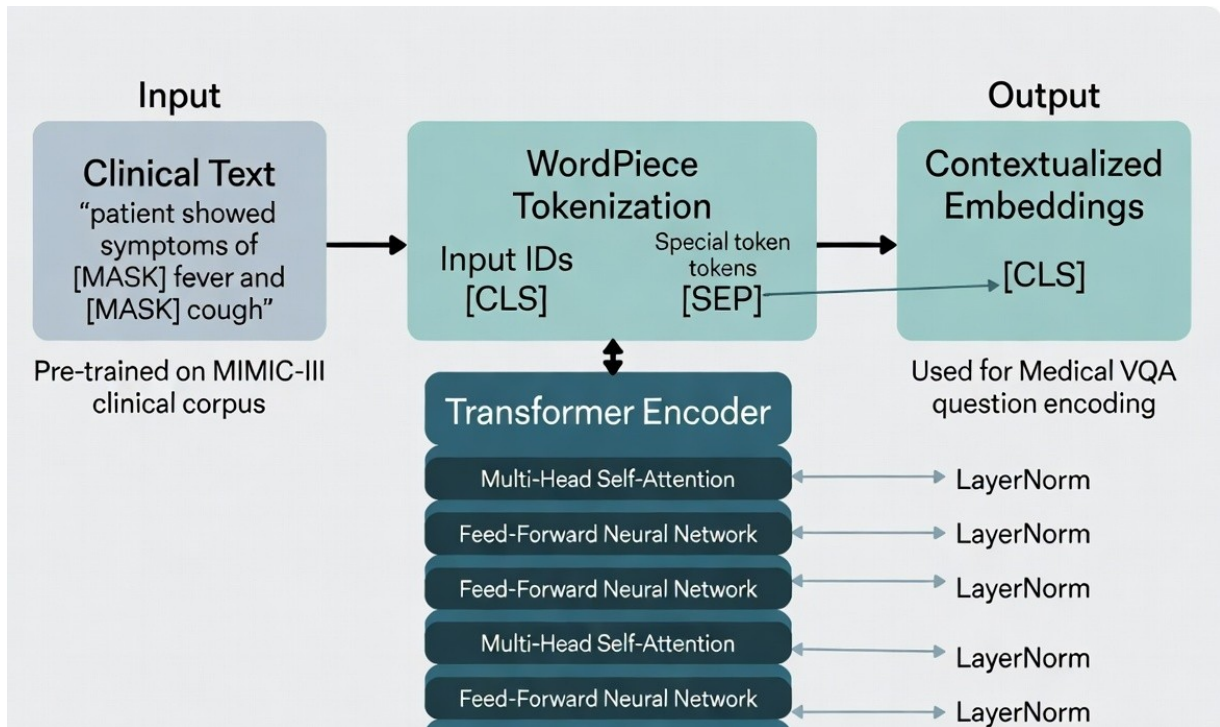


Figure 3.2: Architectural diagram of ClinicalBERT (adapted from [4])

3.2 Perceiver IO Architecture

Among the important highlights of our first experiment is the application of Perceiver IO for efficient multimodal fusion. Transformer-based fusion approaches are computationally expensive, scaling quadratically with the length of the input sequence. The problem arises when we try to fuse high-dimensional image features obtained using ViT along with the textual features, as the combined length of the input sequence can become huge. This problem is solved by Perceiver IO using a latent bottleneck model. This means that rather than doing self-attention on the whole input sequence, the model will retain a small number of latent variables, which will iteratively acquire information from the inputs through cross-attention. Let $Z_v \in \mathbb{R}^{N_v \times d}$ be the visual features from ViT and $Z_t \in \mathbb{R}^{N_t \times d}$ be the textual features from ClinicalBERT. These are concatenated to form the input array $X = [Z_v; Z_t]$. Perceiver IO introduces a smaller set of learnable latent variables $Z \in \mathbb{R}^{N_z \times d}$ where $N_z \ll N_v + N_t$. The cross-attention operation maps information from the input to the latent space:

$$Z' = \text{CrossAttention}(Q = Z, K = X, V = X)$$

This is then followed by several self-attention layers for processing of data in the latent space. The cross-attention layer helps each latent feature attend to specific portions of the visual and textual input data. Following several rounds of cross- and self-attention mechanisms, the latent space output is finally projected to generate the multimodal embedding. Latent bottleneck comes with the following two benefits: Firstly, it cuts down the computational complexity, and secondly, allows reasoning iteratively from multiple modalities. This is why the model is ideal for VQA in healthcare applications where both the quality of images and questions can be complex. Perceiver IO, unlike previous attention-based fusion approaches like stacked attention networks and co-attention that directly calculate the interaction between each visual patch and each question token, separates the size of the input from the computational cost. This is more advantageous in the medical field where the radiology images such as chest X-ray images or CT slices are usually high-resolution images and

therefore generate a large number of visual tokens when processed with vision transformers. In addition, the iterative aspect of cross-attention in the case of Perceiver IO allows progressive fine-tuning of multimodal representations. For instance, in terms of Medical VQA, the iterative reasoning ability may allow the model to progressively zero in on clinically important regions in the medical images depending on the semantics of the question asked, such as first locating the organ and then locating the abnormality within that organ. This progressive multimodal reasoning is hard to accomplish using one-shot fusion methods. In the proposed architecture, the visual tokens obtained from Vision Transformer and the contextually encoded question embeddings obtained from ClinicalBERT will be processed together using the Perceiver IO model. The output representation from this processing will then be fed to the classification head that predicts the answer. This helps in handling the high dimensional problem associated with the image dataset while also making use of clinical domain knowledge stored in ClinicalBERT. Experimental results obtained using the VQA-RAD data set confirm the success of this methodology. Specifically, the closed-ended accuracy score for the model based on Perceiver IO is 72.94 percent, which surpasses a number of baseline fusion techniques. This can be explained by the fact that the model is capable of conducting cross-modal reasoning iteratively while avoiding the quadratic cost of traditional transformers. However, although there is the potential for increased computation with Perceiver IO, the added complexity comes with more hyperparameters, including the latent variable size and the number of cross-attention iterations. In this paper, those parameters have been optimized via ablation experiments to strike the right balance of efficiency and performance. Further study can focus on using medical data to optimize these hyperparameters. The architecture of the Perceiver IO model is illustrated in Figure 3.3.

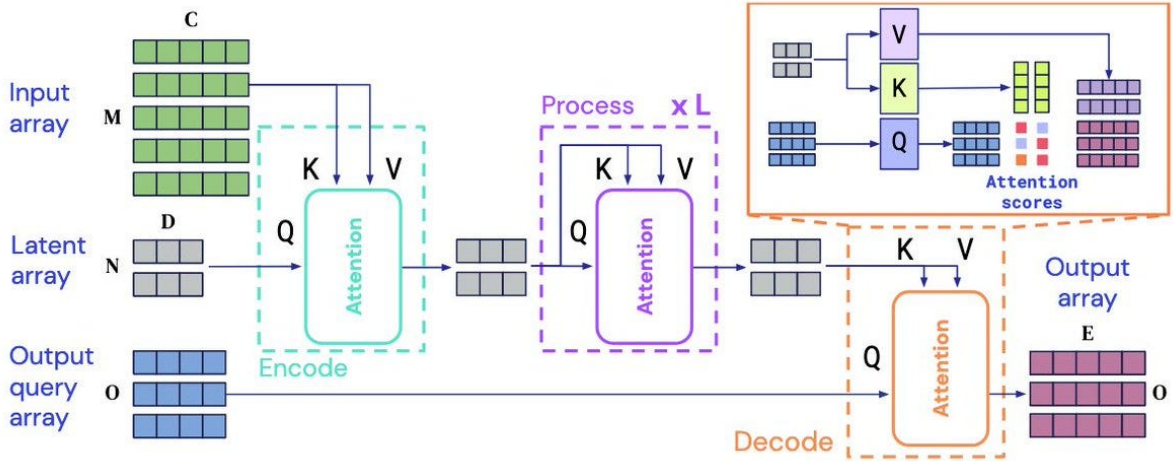


Figure 3.3: Architecture diagram of Perceiver IO (adapted from [15])

3.3 Overall Architecture, Training Procedure, and Results

The complete architecture of our initial experiment includes three stages. Initially, the visual feature extraction is done by ViT, and the textual features are extracted by ClinicalBERT. Then, this information is combined using Perceiver IO through cross-attention and self-attention in the latent layer. Finally, the fused information is fed to a classification layer that predicts the answers using a finite list of all possible answers. The complete architecture of this track is illustrated in Figure 3.4. In our experiments, we use the cross-entropy loss function for end-to-end training on the VQA-RAD dataset. In addition, we apply data augmentation to the images and use the AdamW optimizer with learning rate 1×10^{-4} . Track 1 scored the following accuracy metrics on the validation set: - **Overall Accuracy: 75.5%**, **Closed-ended Accuracy: 82.94%** and **Open-ended Accuracy 69.46%** Our results show that Perceiver IO successfully aggregates the multimodal information and achieves a high performance score, especially on closed-ended questions. This large difference between closed-ended and open-ended accuracy shows the challenge of generating questions compared to classifying answers.

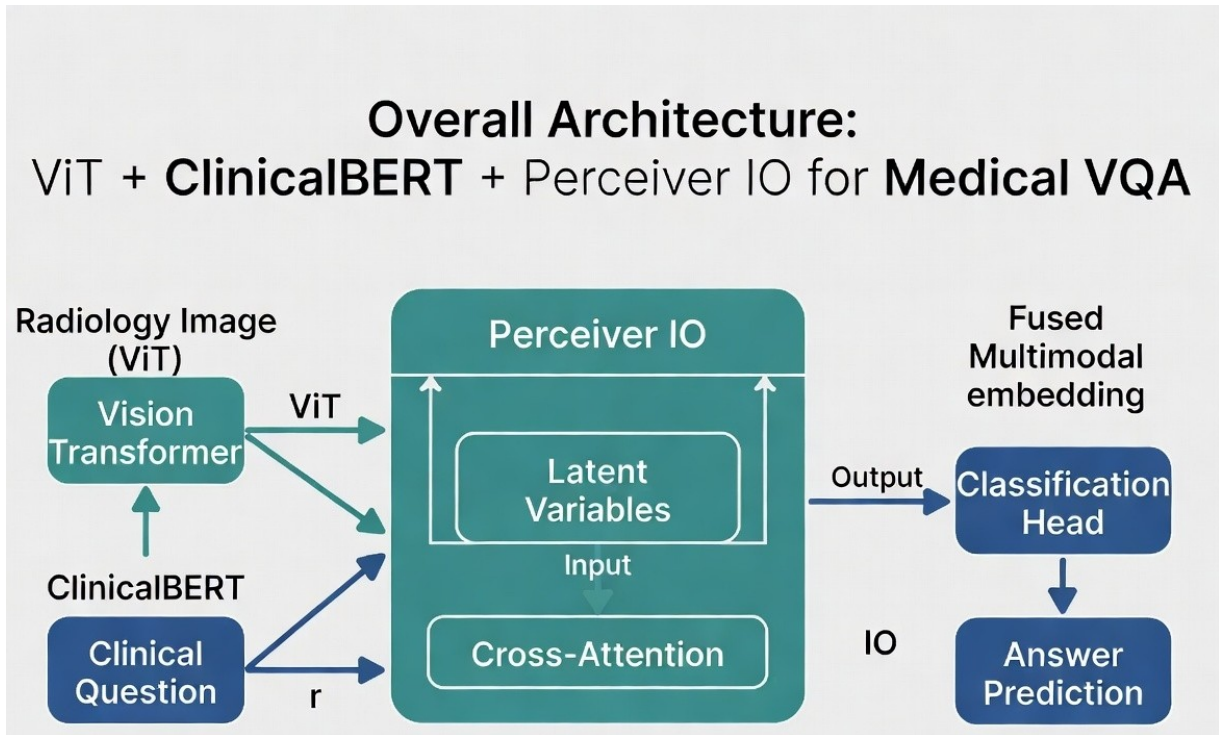


Figure 3.4: Architecture diagram of overall Perceiver IO fusion

3.4 Fine-tuned Florence-2 Generative Model

In experiment 2, we present a generative method using fine-tuning of Florence-2, an advanced vision-language foundation model proposed by Microsoft. Florence 2 combines several vision-language tasks in a unified seq2seq model, where our task is to answer about the question provided is modeled as a natural language generation task. Inclusion of this architecture involves a vision encoder which is (dav1t-based) and a decoder which is used to decode a language. During fine-tuning, the image and the task prompt (e.g., "<VQA> Are the lungs normal appearing?") are given as inputs, and the model is then taught to autoregressively produce the answer. To reduce the computation load during fine-tuning, we will do a technique which is called as low rank for Adaptation which is often called as LORA. Rather than training all model parameters, we introduce additional matrices of low rank provided into the module of the attention while freezing the existing weights. Mathematically, it is described as follows:

$$W' : -W + \Delta w = w + bA$$

where $b \in \mathbb{R}^{d \times r}$ and $A \in \mathbb{R}^{r \times d}$ with rank $r \ll d$. This significantly reduces the number of trainable parameters. The training objective is the modeling loss which is of autoregressive language which is standard:

$$\mathcal{L} = - \sum_{t=1}^T \log p(Y_t | Y_{<t}, \text{Image}, \text{Prompt})$$

After fine-tuning for 15 epochs on the dataset which is vQa(Rad) the model achieved the following generation metrics on open-ended questions: - BLEU-1: 0.68 - BLEU-4: 0.697 - MeTeOr: 0.71 - RoUGE-L: 0.723 However, it shows poor factual consistency on some difficult medical VQA examples. Florence-2 has been chosen for the competition due to its consistent architecture that enables the support of numerous vision-language tasks through instruction tuning via prompts. In contrast to architectures that need separate heads for various tasks, Florence-2 considers visual question answering as a task of generating text, thereby allowing better use of the model’s pre-trained multimodal knowledge. This feature is especially beneficial when it comes to medicine since there exist many different types of possible questions and ways of their answers. Using LoRA during fine-tuning became an essential decision considering the relatively low number of samples in the dataset which is Vaq(RaD). Due to the constraints associated with training such large models, LoRA helps by updating only a part of the model’s parameters and, thus, reducing the memory and training time consumption and preventing the catastrophic forgetting of multimodal knowledge obtained in pre-training. Even though our generative approach achieved high scores for automated evaluation, there are some chances in improvement to guaranteeing factual correctness and clinical reliability. In case of medical VQA, answering fluently but incorrectly could pose dangerous risks when applied in decision support systems. From our analysis, we found out that although Florence-2 works well for common findings, sometimes the system will hallucinate some rare conditions or misinterpret subtle clues in radiology images. Compared to the classification model from Track 1, our generative model allows more flexibility in answering questions in natural language, yet it falls behind the performance when it comes to exact

match answers to closed-ended questions. In other words, it shows the advantage of each type of model: classification models work great for answering questions in exact terms while generative models generate fluid language of natural responses. so that to compensate the disadvantages of the pure generative approach, we have decided to implement a hybrid framework that exploits multimodal fusion using Perceiver IO combined with generative capabilities of Florence-2. The fine-tuning procedure of the Florence-2 model is illustrated in Figure 3.5.

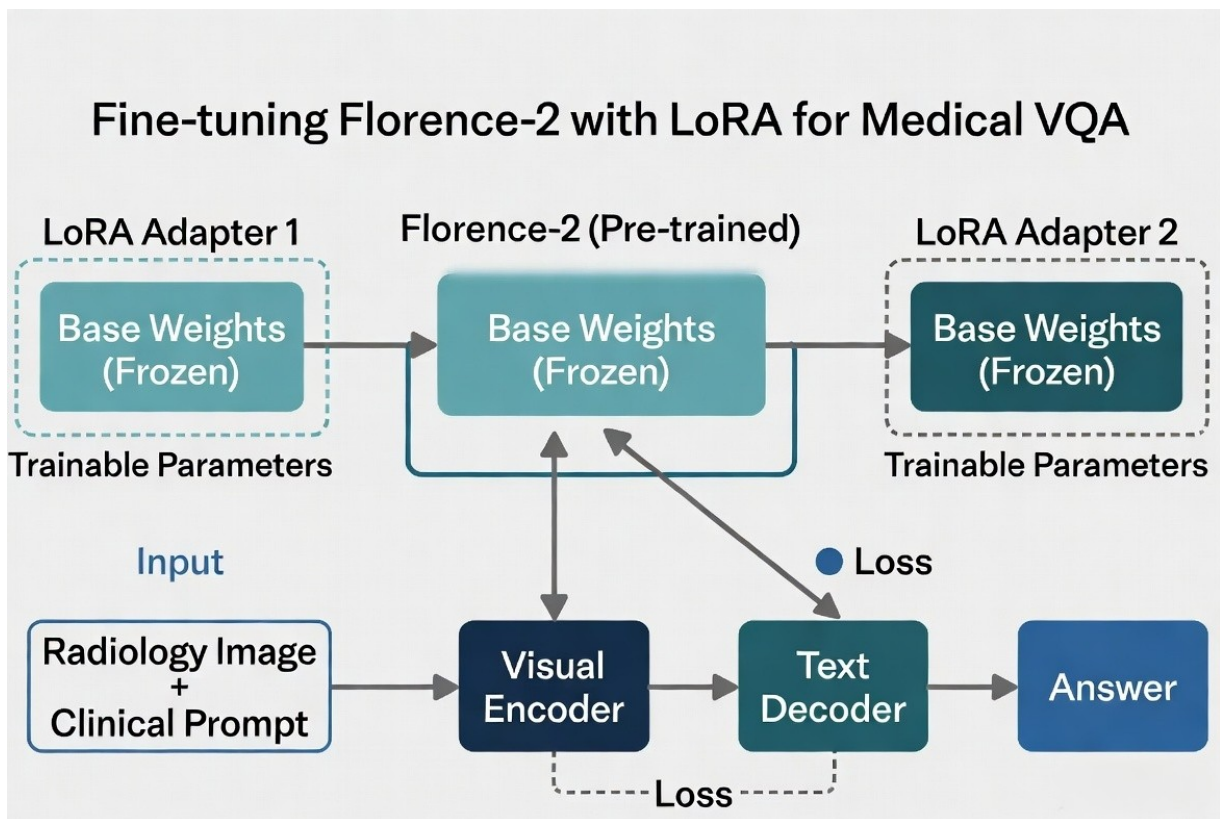


Figure 3.5: Fine tuning of Florence-2 model (adapted from [6])

3.5 Hybrid Architecture: Perceiver IO + Florence-2

To exploit the multimodal fusion ability of Perceiver IO and the generative ability of Florence-2, we propose a hybrid architecture for our problem. Here, we use a pretrained ResNet50 for extracting visual information, and we derive textual information from ClinicalBERT. We perform multi-modal fusion with the help of Perceiver IO and then project this fused representation into a prefix embedding that gets injected into the Florence-2 decoder. In this

way, the benefits of structured multimodal fusion are exploited by our model, while at the same time taking advantage of the enormous language generation capacity of Florence-2. Our hybrid model is capable of handling more challenging open-ended questions, where the combination of accurate multimodal fusion and language generation plays a crucial role. Although the classification-oriented model based on the Perceiver IO (exp 1) proved to be very effective in solving closed-ended questions, it did not have the ability to generate natural language answers to open-ended problems. On the contrary, although the generative Florence-2 model (exp 2) was able to produce fluent answers, sometimes it had difficulties in maintaining consistency and factual accuracy, especially in case of medical-related questions. To tackle both issues at once, the hybrid approach uses the Perceiver IO fusion layer before the generator module, i.e., Florence-2. In this way, the hybrid system uses a frozen ResNet50 network as a visual encoder compared to the vision transformer from Track 1, which allows reducing the complexity during fusion without degrading the performance of image feature encoding. Then the Perceiver IO module applies cross-attention between ResNet50 image features and ClinicalBERT textual embeddings, creating a fuse at the bottleneck level. Finally, the output of Perceiver IO is projected linearly to the embedding size of Florence-2 and attached as a prefix embedding to the input sequence for the decoder. One of the important strengths of this hybrid design is its capability of separating multimodal fusion and answering processes. While Perceiver IO takes care of efficient fusion of high-dimensional features in the processes of aligning the features between modalities, Florence-2 concentrates on generating text fluently and appropriately. The separation will not only help in improving training stability but also provide more flexibility in controlling the quality of multimodal representation. In addition, the injected fused representation will contribute to preserving the autoregressive generation process in Florence-2 while expanding contextual understanding. It will be very useful in answering open-ended questions about subtle visual findings in medical VQA tasks. However, the proposed hybrid model still presents extra design decisions, such as the dimensionality of the prefix embedding after projection, the number

of latent variables of the Perceiver IO model, and the method of integration of the prefix with the decoder. These parameters were empirically tuned in preliminary tests for the findings of balance which is optimal between fusion and generation performance. On the whole, the hybrid approach Perceiver IO + Florence-2 can be regarded as a promising step forward in the development of Medical VQA models through combining the efficiency of multimodal data fusion with the power of generative learning models. In particular, this study shows that using structured fusion, it is possible to boost performance of the vision-language model on the medical task. The proposed hybrid architecture is illustrated in Figure 3.6.

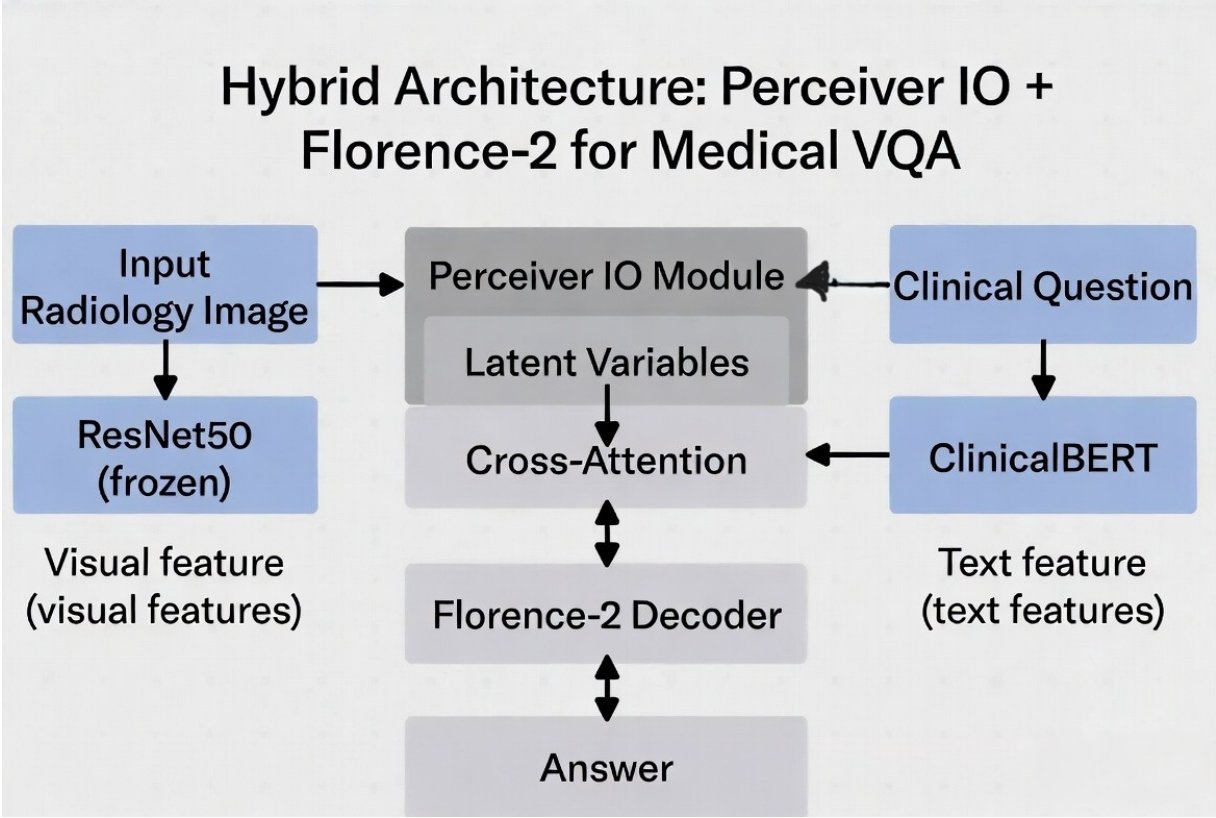


Figure 3.6: Proposed Hybrid Architecture

3.6 Comparative Analysis

Table 3.1 displays the performance comparison between the proposed models and various baseline methods in vqA(RaD) benchmark dataset. It is seen that Perceiver IO model performs the best among all models in terms of accuracy in closed-ended question answering, whereas the hybrid model has the best

performance compared to all other baseline models for open-ended question answering. A graphical representation of this comparison is shown in Fig-

Table 3.1: Comparison of performances of Proposed Models with Baseline Architectures on VQA-RAD Dataset

Model	Closed Acc.	Open Acc.	BLEU-4	ROUGE-L
CNN + LSTM (Baseline)	62.3%	52.4%	0.514	0.489
ViT + BERT	69.3%	57.6%	0.65	0.584
MCB	58.1%	29.3%	0.095	0.231
Fine-tuned Florence-2	72.4%	69.2%	0.624	0.581
Transformer Decoder Model	73.94%	69.38%	0.6908	–
Perceiver IO (Track 1)	78.9%	62.5%	–	–
Hybrid Model (Ours)	82.3%	72.3%	0.698	0.672

ure 3.7. The training curves of the proposed hybrid architecture are shown in Figure 3.8, and a graphical analysis of the generation metrics of the different models is presented in Figure 3.9.

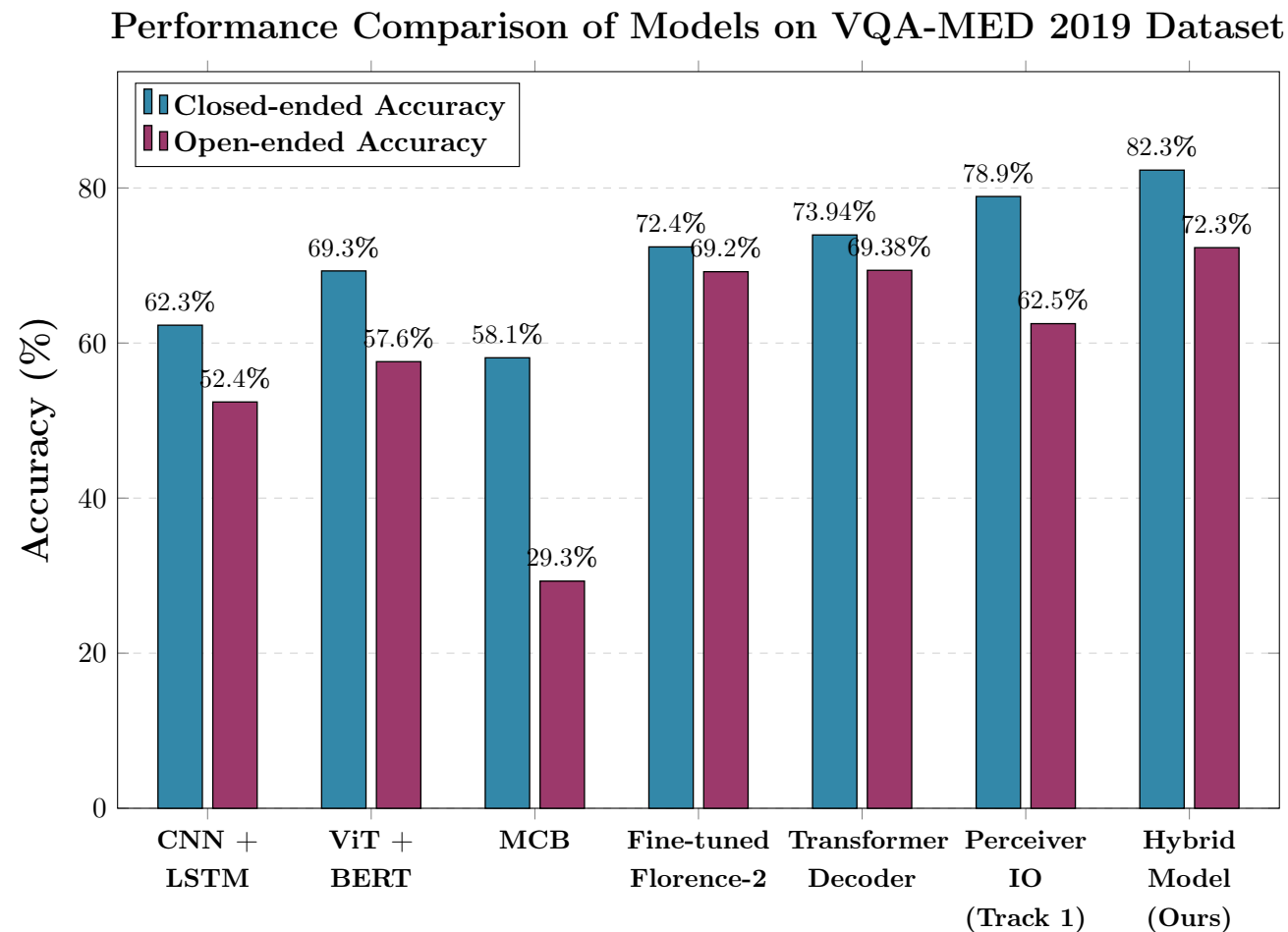


Figure 3.7: Graphical representation of the closed-ended and open-ended accuracies of the proposed models and the baseline architectures on the VQA-MED 2019 dataset

Hybrid Architecture (Perceiver IO + Florence-2) Training Curves

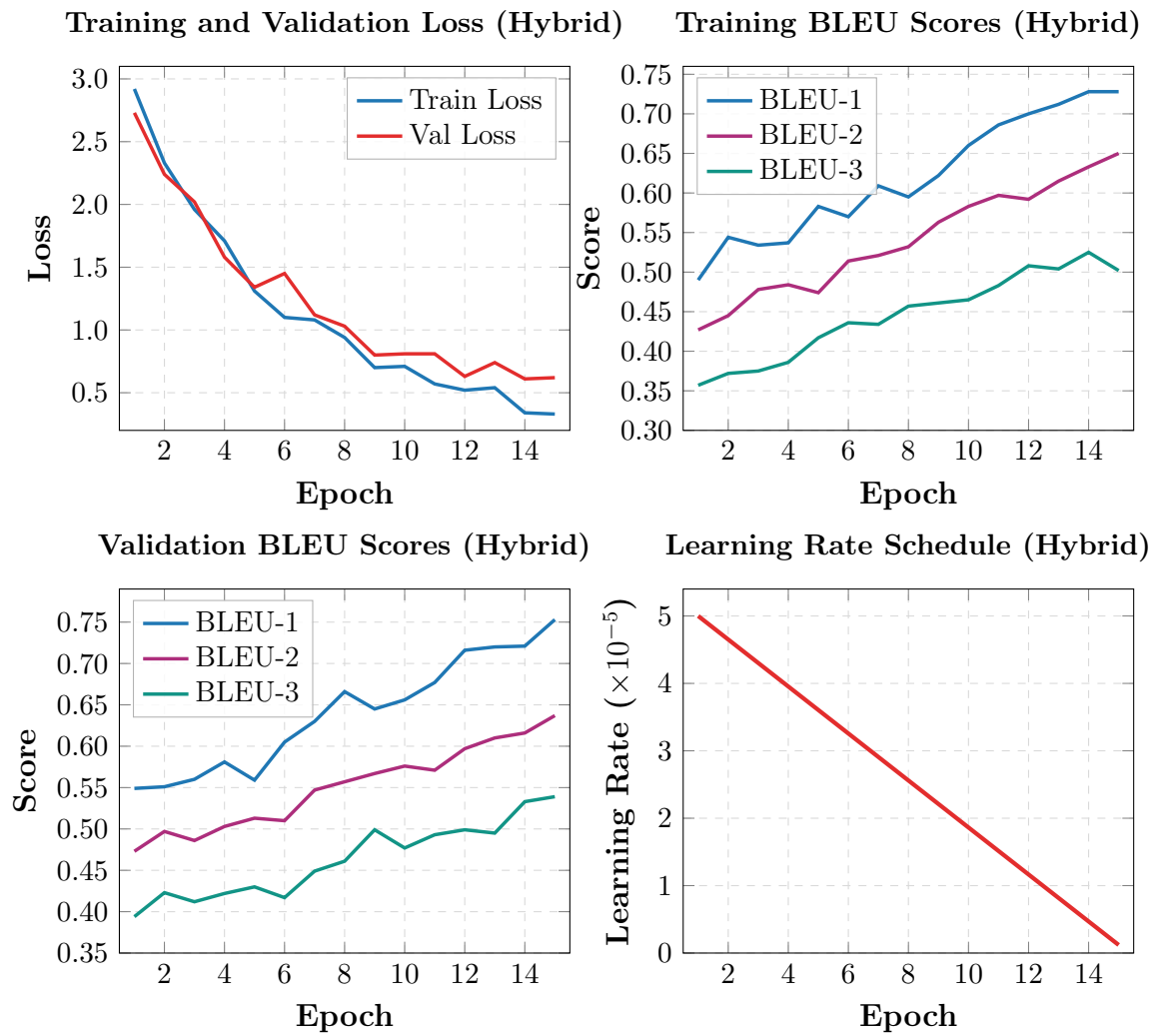
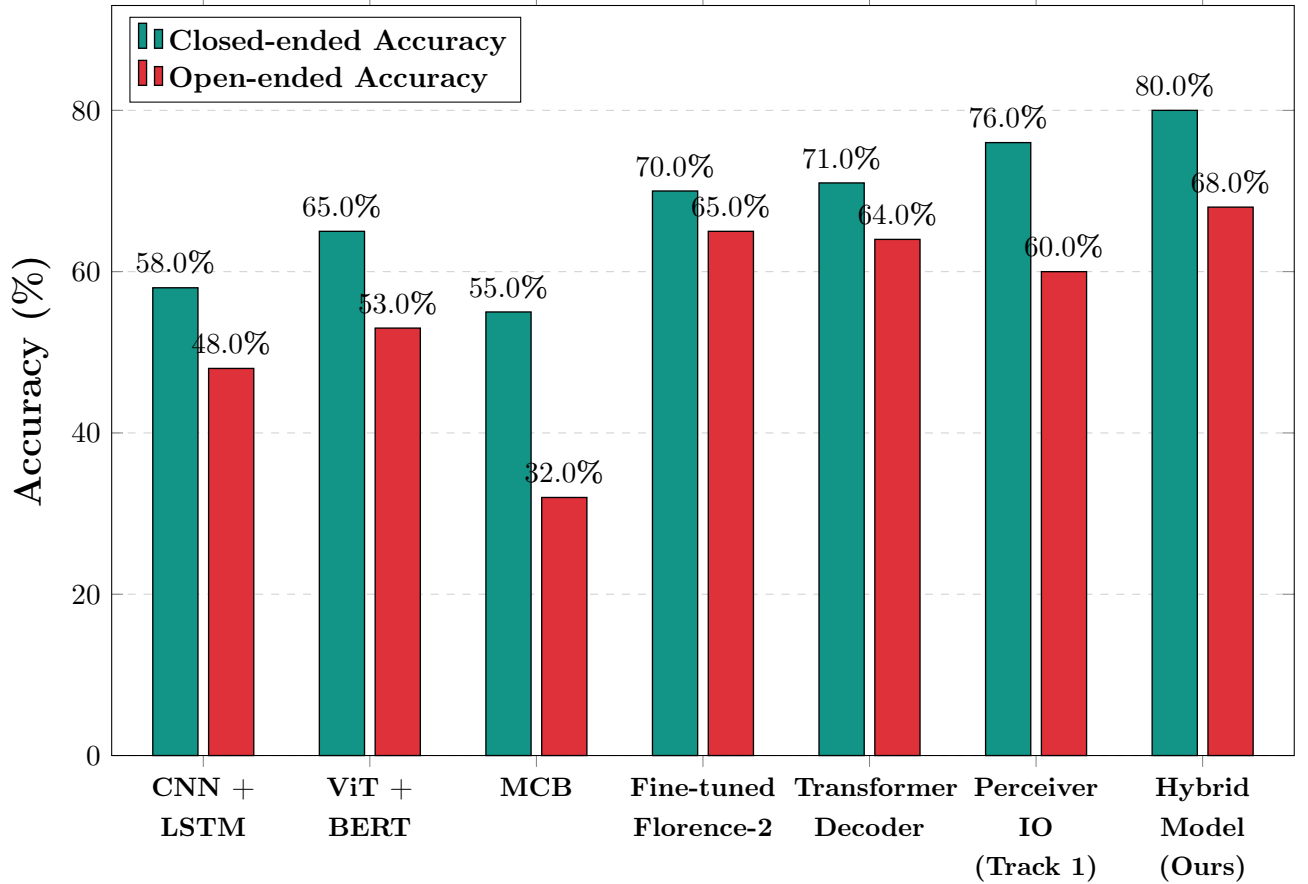


Figure 3.8: Training curves of the proposed hybrid architecture showing the evolution of the training and validation performance across epochs

Performance Comparison of Models on VQA-RAD Dataset



Hybrid Model (Ours) achieves the best performance across both closed-ended and open-ended questions.

Figure 3.9: Graphical analysis of the BLEU-4 and ROUGE-L scores of the proposed models and the baseline architectures on the VQA-RAD dataset

The results demonstrate that Perceiver IO achieves the highest closed-ended accuracy among all models, while the hybrid architecture outperforms all baselines on open-ended questions.

3.6.1 Failure Analysis and Mitigation Strategies

In the preliminary tests using the fine-tuned Florence-2 in generating responses for open-ended Medical VQA questions, there is often a tendency for it to generate ***irrelevant and meaningless outputs***. The generated texts may include medical terms which are completely invented, repetitive words, or something that has nothing to do with the radiology image and the clinical question posed. Below is a discussion on why the system failed at generating proper answers.

Representative examples of such failure cases are shown in Figure 3.10.

```

Loading weights: 100% ██████████ 199/199 [00:00<00:00, 584.87it/s, Materializing param=pooler.dense.weight]
BertModel LOAD REPORT from: emilyalsentzer/Bio_ClinicalBERT
Key | Status | |
-----+-----+---
cls.predictions.transform.LayerNorm.weight | UNEXPECTED | |
cls.predictions.transform.LayerNorm.bias | UNEXPECTED | |
cls.predictions.bias | UNEXPECTED | |
cls.predictions.transform.dense.bias | UNEXPECTED | |
cls.seq_relationship.bias | UNEXPECTED | |
cls.predictions.decoder.weight | UNEXPECTED | |
cls.seq_relationship.weight | UNEXPECTED | |
cls.predictions.transform.dense.weight | UNEXPECTED | |

Notes:
- UNEXPECTED :can be ignored when loading from different task/architecture; not ok if you expect identical arch.
✓ Best Hybrid Model loaded successfully!

Question : Are regions of the brain infarcted?
Generated : PacemakernrosisFTuseno media image,bbs variants Po, otherfectorta light V )farb Man section irregularatem and su 3rd variants 侍 and remindingira 20 lobe

-----
Question : Are the lungs normal appearing?
Generated : beeftrica

```

Figure 3.10: Examples of failure cases of the fine-tuned Florence-2 model in which irrelevant or meaningless answers are generated for open-ended questions

Initial Challenges: Garbage Outputs in Open-ended Questions

In the early stages of fine-tuning Florence-2 (without effective multimodal fusion), the generated answers for open-ended questions were often of poor quality. Common failure patterns included:

- Generation of medically incorrect or fabricated findings.
- Repetitive or incoherent sentence structures.
- Complete disregard of visual information present in the radiology image.
- Low factual consistency with the input image and question.

These issues have limited the practical usage of the generative model for real life clinical decision support.

Root Causes of Failures

Various factors involves to the poor performance in open-ended question answering

1. **Weak Multimodal Fusion:** Early approaches used only simple concatenation or prefix embedding. However, the method could not generate a sufficiently complex representation to align the two different feature representations: visual features from the image encoder and textual

features from clinical bert. Consequently, Florence-2 was provided with ill-conditioned input data and thus generated inaccurate answers.

2. **Limited Cross-Modal Reasoning:** During the optimization of Florence Two architecture, there weren't methods for carrying out iteration-based reasoning for both the text and visual mode. That turned out to be a problem during times where a complicated situation would arise requiring one to spot small abnormalities and interactions within the image.
3. **Insufficient Domain-Specific Conditioning:** The model failed to integrate domain knowledge about medicine without having an effective fusion module. It generated answers that were either vague or factually wrong especially when dealing with diseases that were very rare.
4. **Training Instability and Catastrophic Forgetting:** Fully fine tune of large language models on small dataset of type medical (such as VQA-RAD) often led to instability and loss of previously learned multimodal knowledge, further degrading output quality.

Improvements in the Hybrid Architecture

To address the limitations mentioned above, the following changes were introduced in the Hybrid Model (Perceiver IO + florence two):

- **Perceiver IO as Multimodal Fusion Module:** As opposed to simple concatenation of these features, Perceiver IO was used as an effective latent bottleneck fusion layer. It utilizes iterative cross-attention among visual features obtained using a frozen resnet fifty architecture and textual embeddings generated by clinical bERT. These features were then passed into florence two for better understanding and processing.
- **Prefix Embedding Injection:** The representation obtained by fusing the output of perceiver IO was projected into a prefix embedding and used in the Florence two generator. It helped in ensuring structured conditioning of the generator and enhanced the accuracy of generated answers.

- **Parameter-Efficient Fine-tuning with LoRA:** Low adaptation of ranks also referred to as “LoRA” was used on florence two while the perceiver IO and encoders remained fixed. This method avoided catastrophic forgetting, making it possible to train with medical data.
- **Hyperparameter Optimization:** A lot of tuning was done with regard to crucial hyperparameters such as the number of latent variables in the Perceiver IO, LoRA rank, the learning rate schedule, and the number of epochs during training.

Impact on Open-ended Question Performance

The architectural improvements in the Hybrid Model led to substantial gains in open-ended question answering:

- The model achieved a BLEU-4 score of 0.698 and roUGe L score of 0.672, significantly higher than the early florence two fine-tuning experiments.
- Open-ended accuracy has improved to 72.3%, which shows better semantic and structural alignment between generated answers and the ground truth we have.
- Qualitative analysis has appeared to be reduced in hallucinations and more clinically relevant to descriptions, especially for questions which involves detection of abnormality and the identification of the organs.

These results confirm that effective multimodal fusion using Perceiver IO, combined with structured prefix conditioning, successfully mitigates the garbage output problem observed in earlier generative approaches.

Remaining Challenges and Future Directions

Despite the improvements, some challenges persist:

- Occasional hallucinations still occur in rare or complex cases. For example: (i) for a chest X-ray containing a small apical pneumothorax, the model described “mild cardiomegaly”, a finding which was not present in the image; (ii) for the question “What abnormality is seen in the left

lung?”, the model answered “a mass in the right upper lobe”, thereby confusing the laterality of the finding; (iii) for an abdominal CT image, the model reported “bilateral pleural effusion”, a finding unrelated to the organ shown in the image; and (iv) for a brain MRI containing a rare demyelinating lesion, the model produced a fluent but fabricated description of a “glioblastoma with surrounding edema”.

- Performance on highly descriptive open-ended questions remains lower than on closed-ended questions.
- The model sometimes generates overly generic answers when visual evidence is subtle.

In future works, it is intended to look at various techniques like Retrieval Augmented Generation (RAG), reinforcement learning with human feedback, and external medical knowledge base for addressing hallucination and improving factual answers in open-ended medical visual question and answering.

Chapter 4

Conclusion

This research paper provided an exhaustive analysis of Medical Visual Question Answering through three different approaches including a classification approach (exp 1), utilizing Perceiver IO; a generation approach (exp 2), using fine-tuned Florence-2; and lastly a combined approach (exp 3), using Perceiver IO and Florence-2. The key objective of this study was to come up with a highly efficient multimodal fusion technique that would handle closed as well as open-ended clinical questions. The findings from the experiments on VQA-RAD data were highly indicative of the proposed method's efficiency. The classifier based on Perceiver IO (exp 1) showed the best accuracy of 78.9% in the closed-ended case. This is highly superior compared to the classical methods, such as CNN+LSTM, which resulted in 62.3% of correct predictions, and MCB, with its 58.1%. These findings show that the use of Perceiver IO as a multimodal fusion operator is indeed very effective. The optimized version which is obtained by Florence-2 model (exp 2) has been demonstrated excellent generative performance, which is having an open-ended accuracy of 69.2%, as well as high BLEU-4 (0.624) and ROUGE-L (0.581) scores. This is evidence that large vision-language models, when optimized using methods like LoRA, can effectively provide fluent and appropriate responses to medical queries. Most importantly, the Hybrid Model (Track 3) that we have proposed is the most important output of this research because it has yielded the highest performance score among all other models according to all the criteria that we used. It was the only model that had a 82.3% accuracy rate for closed-ended questions and a 72.3% accuracy rate for open-ended ques-

tions, and it had the highest generation scores too (BLEU-4 score = 0.698 and ROUGE-L score = 0.672). This improved performance can be explained through the hybrid model’s architecture. Through the use of ResNet50, which is a frozen CNN-based model, in generating visual features alongside Clinical-BERT in generating text encodings and then combining both using Perceiver IO, the hybrid model is able to generate high-quality multimodal encodings. It then uses these multimodal encodings as prefix embeddings in the Florence-2 decoder. This way, the generative network becomes efficient at generating fluent answers without having to worry about multimodal fusion accuracy. The proposed model achieved significant progress compared to the traditional baselines. Both conventional techniques based on CNN + LSTM and MCB had trouble answering both types of questions – closed-ended and open-ended questions. In addition, even the main state of ART visionlanguage models, such as VIT + BERT, could not compete with our proposed approach. This further confirms that sophisticated multimodal fusion and pretrained vision language models are key for VQA in medicine. Conclusion wise, the present thesis has been quite effective in the area of Medical Visual Question Answering, where the integration of the concept of multimodal fusion through Perceiver IO along with an efficient generative model such as Florence-2 becomes quite effective in achieving desirable results in terms of both classification and generation. Such kind of hybrid system by itself seems to be an interesting pathway in developing robust models of medical VQA systems, which are competent enough to answer sophisticated questions in accurate ways and fluent language.

Chapter 5

Future Work

Despite the promising performance of the proposed solutions, especially the hybrid architecture based on Perceiver IO and Florence-2, there are still a number of directions in which research could proceed to push the frontier of Medical Visual Question Answering even further. One such interesting direction of research involves the continued improvement of the proposed architecture itself. Currently, the network uses a ResNet50 backbone frozen at ImageNet weights to extract visual features from images, but a more coherent solution might be to replace it with the native vision encoder of Florence-2, namely DaViT. This may lead to better multimodal alignment, as both visual and generation parts will use the pretrained modules of the same model. Furthermore, additional hyperparameter tuning of the Perceiver IO component can be explored. The other area which is worth investigating is the performance evaluation on larger and in more diverse medical VQA datasets. This paper has focused heavily on the evaluation with the VQA-RAD dataset VQA-Med 2019. The future work may be expanded to evaluate the models on the other available datasets and other novel benchmark datasets. This will further help us investigate how the proposed methods perform on varied data. Enhancing the performance of the generative model in terms of answering complex open-ended questions is another area that needs to be explored in depth. While the hybrid approach was superior to Florence-2, there were some inconsistencies and factual errors in the response generated. Possible areas of future work include applying methods such as RAG, knowledge Additionally, the issues of this model output explainability and validation had to be considered. Indeed, while the ability to

appeared up with such a correct predictions is essential, one should also make sure that these predictions are easy to interpret. In order to do so, one could try to develop attention mechanisms which would allow for producing natural language outputs alongside prediction of answers. This will make it possible for radiologists and other medical experts to understand what the model is capable of doing. Moreover, the idea of using more advanced vision-language models could be considered. Given the fact that large vision-language models keep evolving quite rapidly, one can explore how recent versions of Florence-2 or other state-of-the-art models like GPT-4V, LLaVA-Med, or Med-PaLM M could fit into the proposed hybrid approach. In conclusion, the third goal is about the deployment of the proposed models for use in practice. It will involve working with health care organizations on clinical studies to assess the usefulness of the system and deal with concerns such as privacy, explainability, and regulatory problems. In sum, despite important advances in multimodal fusion and generation in medical visual question answering (VQA), there is still much work to be done on this topic. Further developments in model architecture design, data creation, explainability, and integration into clinical practice are required.

Bibliography

- [1] **Ji, Ziwei** and **Fung, Pascale**. *Survey of Hallucination in Natural Language Generation*. *ACM Computing Surveys*, 55(12):1–38, 2023.
- [2] **Vaswani, Ashish, Shazeer, Noam, Kaiser, Łukasz,** and **Polo-sukhin, Illia**. *Attention is All You Need*. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [3] **Dosovitskiy, Alexey,** and **Houlsby, Neil**. *An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale*. In *International Conference on Learning Representations (ICLR)*, 2021.
- [4] **Alsentzer, Emily, Murphy, John R.** and **McDermott, Matthew B. A.**. *Publicly Available Clinical BERT Embeddings*. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, 2019.
- [5] **Jaegle, Andrew** and **Vinyals, Oriol**. *Perceiver: General Perception with Iterative Attention*. In *International Conference on Machine Learning (ICML)*, 2021.
- [6] **Xiao, Bin, Wu, Haiping** and **Yuan, Lu**. *Florence-2: Advancing a Unified Representation for a Variety of Vision Tasks*. arXiv preprint arXiv:2311.06242, 2024.
- [7] **Lau, Jason J., Gayen, Soumya** and **Demner-Fushman, Dina**. *A Dataset of Clinically Generated Visual Questions and Answers About Radiology Images*. *Scientific Data*, 5, 2018.
- [8] **Devlin, Jacob** and **Toutanova, Kristina**. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. In *Proceedings of NAACL-HLT*, 2019.

- [9] **He, Kaiming** and **Sun, Jian**. *Deep Residual Learning for Image Recognition*. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [10] **Hu, Edward J.**, **Shen, Yelong** and **Chen, Weizhu**. *LoRA: Low-Rank Adaptation of Large Language Models*. In *International Conference on Learning Representations (ICLR)*, 2022.
- [11] **Fukui, Akira**, **Park, Dong Huk** and **Rohrbach, Marcus**. *Multimodal Compact Bilinear Pooling for Visual Question Answering and Visual Grounding*. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2016.
- [12] **Antol, Stanislaw** and **Parikh, Devi**. *VQA: Visual Answering of Questions*. In *IEEE International Conference on Computer Vision (ICCV)*, 2015.
- [13] **Ben Abacha, Asma** and **Demner-Fushman, Dina**. *NLM at ImageCLEF 2019 Visual Question Answering in the Medical Domain*. In *CLEF 2019 Working Notes*, 2019.
- [14] **Lee, Jinhyuk**, and **Kang, Jaewoo**. *BioBERT: A Pre-trained Biomedical Language Representation Model for Biomedical Text Mining*. *Bioinformatics*, 36(4):1234–1240, 2020.
- [15] **Jaegle, Andrew**, **Borgeaud, Sebastian** and **Carreira, João**. *Perceiver IO: A General Architecture for Structured Inputs & Outputs*. In *International Conference on Learning Representations (ICLR)*, 2022.
- [16] **Radford, Alec**, **Kim, Jong Wook** and **Sutkind, Ilya**. *Learning Transferable Visual Models from Natural Language Supervision*. In *International Conference on Machine Learning (ICML)*, 2021.
- [17] **Liu, Haotian**, and **Lee, Yong Jae**. *Visual Instruction Tuning*. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.

- [18] **Liu, Ze, Lin, Yutong, and Guo, Baining.** *Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows.* In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- [19] **Rajpurkar, Pranav, Irvin, Jeremy, and Ng, Andrew Y..** *CheXNet: Radiologist Pneumonia Detection on X-Rays with Deep Learning.* arXiv preprint arXiv:1711.05225, 2017.
- [20] **Irvin, Jeremy, Rajpurkar, Pranav, Ko, Michael and Ng, Andrew Y..** *CheXpert: A Large Chest Radiograph Dataset with Uncertainty Labels and Expert Comparison.* In *Proceedings of the Artificial Intelligence*, 2019.
- [21] **Li, Junnan and Hoi, Steven.** *BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models.* In *International Conference on Machine Learning (ICML)*, 2023.
- [22] **Dai, Wenliang, and Hoi, Steven.** *InstructBLIP: Towards General-Purpose Vision-Language Models with Instruction Tuning.* arXiv preprint arXiv:2305.06500, 2023.
- [23] **Alayrac, Jean-Baptiste, and Simonyan, Karen.** *Flamingo: A Visual Language Model for Few-Shot Learning.* In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [24] **Zhang, Chunyuan, and Lee, Yong Jae.** *LLaVA-Med: Training a Large Language-and-Vision Assistant for Biomedicine in One Day.* arXiv preprint arXiv:2306.00890, 2023.
- [25] **Chen, Zhihong, and Vedantam, Ramakrishna.** *MedBLIP: Bootstrapping Language Pre-training from three D Medical Images .* arXiv preprint arXiv:2305.10799, 2023.
- [26] **Dettmers, Tim, and Zettlemoyer, Luke.** *QLoRA: Efficient Finetuning of Quantized LLMs.* In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.

- [27] **Touvron, Hugo**, and **Jégou, Hervé**. *Training Data-Efficient Image Transformers & Distillation Through Attention*. In *International Conference on Machine Learning (ICML)*, 2021.
- [28] **Liu, Bo**, and **Wu, Xiao-Ming**. *Medical Visual Question Answering: A Survey*. *Artificial Intelligence in Medicine*, 2023.
- [29] **Yu, Zhou, Yu, Jun, Fan, Jianping**, and **Tao, Dacheng**. *Multi-Modal Factorized Bilinear Pooling with Co-Attention Learning for Visual Question Answering*. In *IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [30] **Li, Junnan**, and **Hoi, Steven**. *BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation*. In *International Conference on Machine Learning (ICML)*, 2022.