# COMPUTER-AIDED MEDICAL DIAGNOSTIC SYSTEM ON VISUAL QUESTION-ANSWER USING DEEP LEARNING

by

Nihar Ranjan Behera, DBA Research Scholar

DISSERTATION

Presented to the Swiss School of Business and Management

Geneva in Partial Fulfillment

Of the

Requirements

for the

Degree

DOCTOR OF BUSINESS ADMINISTRATION
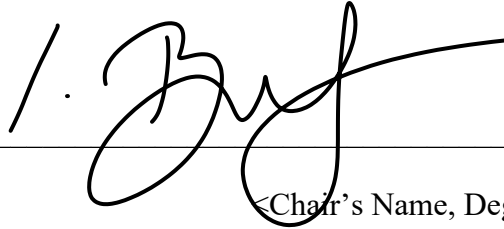
SWISS SCHOOL OF BUSINESS AND MANAGEMENT

GENEVA May 2023

# COMPUTER-AIDED MEDICAL DIAGNOSTIC SYSTEM ON VISUAL

# QUESTION-ANSER USING DEEP LEARNING

by

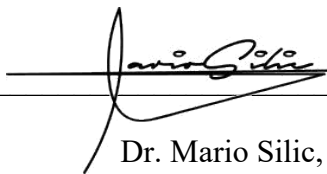Nihar Ranjan Behera, DBA Research Scholar

APPROVED BY

_____

<Chair's Name, Degree>, Chair

_____

<Member's Name, Degree>, Committee Member

_____

Dr. Mario Silic, Ph.D., Research Supervisor

RECEIVED/APPROVED BY:

_____

<Associate Dean's Name, Degree>, Associate Dean

2

## Dedication

I would like to dedicate this work to my mother late Anusuya Behera, who passed away recently, she is always a source of light in the darkest and toughest days of my life. I am grateful to learn the lessons of determination, honesty, and compassion she taught me that helps to grow.

**DECLARATION**

First and foremost, I want to thank my company, "Innovasoft consultants Limited, Ireland" for allowing me to work for many years in the artificial intelligence and machine learning domain in many different countries, where I have learned about the practical issues that various organizations face when implementing artificial intelligence and machine learning. I also want to express my gratitude to my coworkers, clients, and mentors for their ongoing encouragement and support during this journey. Second, I want to express my gratitude to Professor Dr. Mario Silic, Ph.D., the dissertation chair of my research and supervisor, for his unwavering support, counsel, criticism, and encouragement during the production of my dissertation. Thirdly, I'd want to take this chance to express my gratitude to the SSBM administration and staff for giving me the chance to attend a top-tier business school and to all the doctorate students and colleagues for their assistance during the research. I want to express my gratitude to all the businesses that provided this research study with their invaluable contributions. I would not have been able to understand actual issues and further pinpoint remedies without their genuine contribution of essential Artificial Intelligence, and Machine Learning material on their own. Last but not least, I want to mention how much my wife Alka, my children Aditi and Advik, as well as my father Mr. Muralidhar Behera, supported me when I was studying after office hours and letting to concentrate on my research studies while pulling quite a significant amount of quality family time. I want to dedicate this research to all of them.

**ABSTRACT**

COMPUTER-AIDED MEDICAL DIAGNOSTIC SYSTEM
ON VISUAL QUESTION-ANSER USING DEEP LEARNING

by

Nihar Ranjan Behera

2023

Dissertation Chair:

Co-Chair:

Medical-VQA can enhance the capabilities of computer-aided diagnosis that can not only bridge the imbalance in the current state of medical resources in the healthcare sector but also reduce the cost for both the hospitals as well as for the patients. Medical images are rapidly being used in medical diagnosis, however, VQA systems are lacking behind due to capabilities of the domain transfer of general VQA model to medical-VQA as there is a huge gap between the nature and complexity of general domain and medical domain datasets. Moreover, medical-VQA datasets are also not very big which also enhances the need for such a VQA system that can learn from the lesser data and can be generalized to diversify medical imaging. In this proposal, a multimodal Transformer based VQA system is suggested. The suggested system incorporates the interaction among input images, question and answer text.

**TABLE OF CONTENTS**

LIST OF TABLES

## LIST OF FIGURES

# CHAPTER 1.

# INTRODUCTION

## 6.1 Introduction

Artificial intelligence (AI) is not only making our lives uncomplicated by automating a tedious and time-consuming task, but it also has various uses in medical health care to diagnose people. One of our primary immediate concerns has always been health. So far, there have been few simple methods for finding out about our physical conditions without expert help. Imaging, as a non-invasive approach to creating images of the internal aspects of the body, is an essential tool for doctors in clinical diagnosis and evaluation. Numerous patients have unanswered questions about surgical procedures and medical diagnoses and the majority of these questions go unanswered because due to the limited number of medical experts who are frequently overburdened with academic and clinical work, their chances of finding a health professional to confirm their doubts are slim (Bates & Gawande, 2000). Many desktop techniques ("Computer-Assisted Surgery," n.d.)(Rogers et al., n.d.) and simulators (Kneebone, 2003; Sarker et al., 2007) have been proposed to help students improve their surgical skills.

Even though the systems help students to improve their abilities and minimize the workloads of academic professionals, they do not attempt to respond to students' questions. Whereas students have already been known to learn by observing recorded surgical techniques, it is still the responsibility of medical experts to answer their

questions. In such situations, a computer-assisted system that can perform both questionnaires and medical data and provide accurate answers that would benefit the pre-screening process greatly while minimizing the workload of the medical expert (Seenivasan et al., 2022).

Proposed medical-VQA systems (Zhan *et al.*, 2020; Abacha *et al.*, 2018; Eslami, de Melo, and Meinel, 2021) are based on transfer learning (Pan, data and 2009) and exploiting the existing Deep learning based VQA system by fine-tuning the pre-trained models on the comparatively smaller medical dataset that are pre-trained on the general datasets. However, there is a very gap between the general and medical domains, both in terms of image structure as well as the language of question-answer pairs (Eslami, de Melo, and Meinel, 2021). So, the transfer of domain knowledge from the general to the medical domain through transfer learning (Pan, data and 2009, no date) and fine-tuning is not sufficient enough to address the medical questions. In this research, I am to build a Visual question answering for medical images that will answer the asked medical queries, both close-ended and open-ended, through analyzing the radiology images. A trilinear interactive attention-based method proposes in this study to improve each individual modality depiction by combining other modalities (TrI-Att). Furthermore, while self-attention cannot fuse various modalities, it can improve their interaction (Yu et al., 2019). The input image will be supplied into Efficient Net, while the questions and answers will be gone through a pre-trained BERT to create an image and text embedding, respectively. The embedding then is processed using trilinear attention to concatenate multimodal

input. Individual and combined inputs are gone through the residual MLP layers, followed by self-attention layers. The last residual MLP output is fed to the Transformer encoder and then to the output layer to predict the answer. On inference, the proposed VQA receives the question and image and generates an answer.

Surgical scenes are enhanced with data that the system can use to answer questionnaires about defective tissue, surgical tool conversation, and surgical procedures. With the ability to extract various types of data from a simple image feature simply by changing the question, the computer vision field has recently seen an influx of vision and NLP models for medical-VQA systems (L. Li et al., n.d.; Z. Wang et al., n.d.). These models are based on either the long short-term memory (LSTM) (Barra et al., n.d.; Himanshu Sharma & Jalal, 2021)or attention methods (Seenivasan et al., 2022; H Sharma et al., n.d.).

Different doctors may obtain different information from medical images. Deep learning, as an efficient information processing method, is becoming more important in health informatics (Justin Ker; Lipo Wang; Jai Rao; Tchoyoson Lim, 2017). In contrast to the computer vision field, which is frequently supplemented with massive, annotated datasets, the medical field lacks a sufficient amount of annotated data, reducing medical VQA exploration. Because of the unavailability of domain-specific clinical terms, The transfer learning methods alone cannot be sufficient to modify pre-trained computer-vision VQA systems for clinical applications. While some work on medical-VQA

(Seenivasan et al., 2022) for clinical diagnosis has recently been reported, VQA for surgical images remains largely unknown.

In the medical profession, a good deep learning-based VQA model can automatically retrieve information from medical images and aid in clinical diagnosis. Meanwhile, the medical VQA model can assist patients in gaining a basic understanding of their physical state, which can then be used to select a more focused medical treatment strategy. In general, implementing the medical VQA model can help to alleviate the problems created by the imbalanced distribution of healthcare resources. In past years, the use of computer-aided diagnosis (CAD) techniques for processing medical data has grown in popularity (graphics & 2007, n.d.). However, many CADs focus on histopathologic diagnosis or segmentation of a specific type of medical image, including tumor tracking (R. Wang et al., n.d.). Some CADs use medical records to estimate risks (Ren et al., n.d.) but the majority of CAD processes are designed to diagnose a single disease, such as breast cancer (Bardou et al., n.d.), or lung disease (C. Li et al., n.d.). Figure. 1 shows the medical-VQA scheme.

| | Closed-ended | Open-ended |
|---|---|---|
| Question | Are the clavicles broken? | Where is the lesion located? |
| Answer | No | Anterior mediastinum |
| Question Type | Object condition presence | Position |
| Organ | Chest | |

| | Closed-ended | Open-ended |
|---|---|---|
| Question | Is the appendix normal in size? | What cut of the body is this image? |
| Answer | Yes | Axial |
| Question Type | Size | Plane |
| Organ | Abdomen | |

*Figure 1:Close-ended and Open-ended samples of VQA input and output.*

There is, however, less work being carried out to combine NLP with medical image processes, including medical image captioning (Eickhoff et al., n.d.) and medical VQA (Hasan et al., 2018). VQA-RAD (Lau et al., 2018), which includes 315 medical images, is an example. This dataset contains a collection of questions ranging from simple to complex. Tremendous questions concern the severity of the disease or the diagnosis. This dataset contains a wide variety of healthcare images and question-answer pairs representative of the real-world medical environment.

As per recent studies, the medical VQA has been directed to various "jobs." The first is the diagnosing radiologist, who consults with the referring physician as an expert. According to a workload analysis (McDonald et al., n.d.), the average radiologist must perceive one CT or MRI image in 2 to 4 minutes. A radiologist must answer approximately 27 phone calls every day from patients and physicians in addition to the

extensive queue of imaging studies ("The Voice of the Radiologist: Enabling Patients to Speak Directly to Radiologists," n.d.), resulting in additional disruptions and inefficiencies in the workflow. A medical VQA system has the potential to answer physicians' questions, alleviate the burden on the healthcare system, and improve the efficiency of medical professionals.

Another use for VQA is to respond as pathologists, examining body tissues and assisting other healthcare professionals in making diagnoses (He et al., 2020). In addition to the role of a healthcare professional, the medical-VQA system can act like a knowledgeable assistant. For instance, the "second opinion" from the medical-VQA system can help clinicians interpret medical images while also lowering the risk of misdiagnosis. Finally, (Tschandl et al., n.d.) a fully developed and accomplished medical-VQA system can directly evaluate images of patients and answer relevant types of questions. A medical-VQA system can provide equivalent consultation in some instances, such as completely automated health assessments, where medical professionals might not be accessible. Just after a hospital visit, patients look for additional information online. The search engine's misleading and irregular information may lead to wrong answers. A medical VQA can also be embedded through an online consultation process to provide accurate answers at any time and from anywhere.

Medical VQA is more difficult than general-domain VQA due to the following factors. To begin, developing a large healthcare VQA dataset is difficult because professional

annotation is expensive due to the high level of professional knowledge required, and QA sets cannot be created directly from images. Second, answering questions focused on a medical image necessitates a unique design of the VQA model. Because a diagnosis is microscopic, the task must also concentrate on a fine-grained dimension. As a result, segmentation methods may be needed to precisely locate the region of interest. Ultimately, a question can be very skilled, requiring the model to be trained with medical expertise instead of a general language database.

Medical imaging (MI) (Yuan *et al.*, no date) is a domain of producing images of the internal vision of the patient's body and is a very significant tool for clinical diagnosis and disease analysis. MI helps the doctor to visualize and analyze the patient's physical conditions to diagnose and assess the irregularities. However, the analysis and interpretation of different medical practitioners on the same medical images can be different. Deep Learning (Goodfellow, Bengio, and Courville, 2016) based models are very powerful to process complex data like medical images, and can greatly benefit the information processing in health informatics. State-of-the-art deep learning methods are available to extract the desired information from visual or textual data. In the context of medical imaging, a good deep learning-based medical VQA mechanism can extract the desired information from the medical images to answer the asked question and help in the disease diagnosis. Moreover, patients can benefit from the VQA system by getting a preliminary understanding of their body condition which can guide them in choosing a more precise treatment plan. The examples of med-VQA are shown in Table 1.

Visual question-answering (VQA) (Wang *et al.*, no date) has become a very prominent area of research in Deep Learning due to immense development in Natural Language Processing (NLP) and Computer Vision (CV). VQA is an amalgamation of NLP and CV to join to understand the knowledge representation between both domains by transferring and corresponding to the knowledge space. VQA aims to answer natural language queries by extracting and then corresponding to the relevant visual information. VQA techniques are applied to understand the scene setting in the image and then address the asked questions about that scene (Goyal *et al.*, no date). Applying the VQA techniques to daily scenes, object detection, image classification, or attention mechanism can be effective (Antol *et al.*, no date).

In general, an effective medical VQA (Abacha *et al.*, 2019; Abacha *et al.*, no date; Lau *et al.*, 2019) system can bridge the gap of the imbalanced distribution of medical resources. However, medical VQA lacks a giant labeled dataset that is the prerequisite for the training of deep learning models. Another challenging aspect is the medical-related vocabulary that usually contains very different words as compared to the daily life language. So, medical terminologies are difficult to interpret and formulate a meaningful query against which relevant visual information can be extracted to answer the asked question.

Table 1:Examples from the dataset of VQA-Med2019, VQA-Med2020, Path-VQA, and VQA-RAD, respectively data having Images with QA pair.

19

| | | |
|---|---|---|
|  |  |  |
| Q: What modality is shown? A: CTA-CT angiography. | Q: What imaging modality was used to take this image ? A: CT with iv contrast. | Q: What organ system is displayed in this ct scan? A: Skull and contents. |
|  |  |  |
| Q: what is abnormal in the ultrasound? | Q: what abnormality is seen n in the image? A: Aberrant right subclavii -an artery (ARSA). | Q: what is most alarming about this ct scan? A: Medullary nephrocalcinosis. |

| | | |
|---|---|---|
| A: Perihepatic fluid and fluid in the gallbladder as expected for the postoperative patient. | | |
|  |  |  |
| Q: What does the sectioning of the ovary show?<br><br>A: large endometriotic cyst with degenerated blood. | Q: What does the mucosal surface show?<br><br>A: Papillary tumor floating in the lumen. | Q: What are inactive facto -rs?<br><br>A: The red polypeptides. |
|  |  |  |

| Q: In what plane is this image taken? | Q: Which plane is this image taken? | Q: How would you describe the abnormalities |
|---|---|---|
| A: Axial. | A: PA. | ?<br><br>A: Ring-enhancing lesions<br><br>. |

In Table. 2, there are the details of medical image datasets that contain image and QA pair. These are the three most comprehensive radiology VQA dataset that involves multiple diseases and has almost all the most commonly asked types of questions by the medical fraternity. Table. 3 describes each included question type given for a better understanding of the scope of the suggested dataset. For this study, these three datasets can be combined to make a comprehensive dataset. Jointly, there will be 4436 images having around 17, 000 QA pairs. QA pairs include closed and open-ended questions. These images and QA pairs have great diversification in terms of various body organ images, and all the major asked questions by radiologists.

Table 2:Details of the datasets used for the training and evaluation of the proposed medical-VQA

| Dataset | labeled images | No of QA pairs | Question type | Answer types |
|---|---|---|---|---|

| | | | | |
|---|---|---|---|---|
| VQA_RAD (Lau et al., 2018) | 315 | 2248 | MODALITY PLANE ORGAN (Organ Syst-em) ABN (Abnormality) PRES (Object/Condition Presence) POS (Positional Reaso-ning) COLOR SIZE ATTRIB (Attribute Ot-her) COUNT (Counting) Other | Closed-ended Open-ended |
| VQA-Med-2019 (A. Abacha et al., 2019) | 3,200 medical images | 12,792 | Modality Plane Organ system Abnormality | Closed-ended Open-ended |

| VQA-Med-2020 (A. Abacha *et al.*, 2020) | 921 | 2,320 | Modality Plane Organ system Abnormality | Closed-ended Open-ended |
| --- | --- | --- | --- | --- |

*Table 3: Details of the datasets used for the training and evaluation of the proposed medical-VQA*

| Question Type | Description |
| --- | --- |
| Modality | How an image is taken – CT, x-ray, T2 weighted MRI, etc. |
| Plane | Orientation of an image slicing through the body – axial, sagittal, coronal |
| Organ System | The categorization that connects anatomical structures with pathophysiology, diagnosis, and treatment – pulmonary, cardiac, musculoskeletal system |
| Abnormality | The normalcy of an image or object. For example, "is there something wrong with the image?" or "What is abnormal about the lung?", "Does the liver look normal?" |
| Object/Condition Presence | Objects could be normal structures like organs or body parts but could also be abnormal objects such as masses or lesions. Clinicians |

| | may refer to the presence of conditions in an image or patient – fractures, midline shift, infarction |
|---|---|
| Positional reasoning | Position or location of an object or organ, including what side of a patient, in respect to the image borders, or relative to other objects in the image |
| Color | Signal intensity including enhancement or opaqueness |
| Size | Measurement of the size of an object, e.g., enlargement, atrophy |
| Attribute Other | Other types of description questions |
| Counting | Focusing on several objects, e.g., the number of lesions |
| Other | Catch-all categorization for questions that do not fall into the previous categories |

Answering questions based on medical images need to have an understanding of vision processing, natural language processing, and medical knowledge base. To enhance the medical-VQA understanding, some models employ knowledge-based reasoning for VQA. VQA models based on computer vision, such as classification and object detection, are not required to extract and understand the image information fully but are insufficient to answer relatively complex questions and mostly do well on closed-ended questions only. When medical-VQA (Abacha *et al.*, 2019; Abacha *et al.*, no date) compared with the

general daily life images-based VQA system, medical VQAs (Abacha *et al.*, 2019) (Abacha *et al.*, no date) is a much more complex and riskier task because medical questions are difficult to contextualize with images, and also requires to answer the answer very accurately as the health and wellbeing of the patient depend on. Medical-VQA (Abacha *et al.*, 2019; Abacha *et al.*, no date) systems should have the ability to handle the complexity and highly accurate performance simultaneously.

Visual perception-based medical-VQA (Abacha *et al.*, 2019; Abacha *et al.*, no date) tasks such as recognition of modality, localization of particular lesions, or determining the normality of a patient's organ, require extraction of high-level features that can provide the reasoning to address the question. To answer the medical questions accurately, reasoning ability and medical domain knowledge are the prerequisites. Another great problem is the lack of an accurately labeled medical-VQA (Abacha *et al.*, 2019; Abacha *et al.*, no date) dataset for the Deep Learning based model training because labeling the images by a medical expert is difficult due to his inability to understand how the Deep learning model works that's why it's very laborious and expensive to obtain high-quality annotations by medical experts.

Our VQA system is useful for Business medical intelligence, in which the system is used to extract insights from medical images to improve the efficiency and profitability of healthcare organizations, such as identifying trends in patient conditions or detecting patterns in treatment outcomes, also can be used to provide interactive, multimedia-based

26

training and education for medical students and residents, allowing them to learn and practice diagnostic skills using real-world medical images and can be integrated into telemedicine platforms, allowing doctors to remotely diagnose and treat patients using medical images, which can be particularly useful for patients in remote or underserved areas. So, in this way, the system will not only save 100 to 500 dollars for the patient but also examine the patient's radiographic image for each potential medical problem.

This research helps businesses in the following ways:
- o Improved Diagnostic Accuracy
- o Improvements in Efficiency and Productivity
- o Cost Reduction
- o Knowledge Sharing and Collaboration
- o Continuous Learning and Improvement
- o Competitive Advantage
- o Augmented Decision-Making
- o Enhanced Patient Care and Experience

## 6.2 Research Problem

Medical VQA can be very beneficial both for the medical community as well as for the patient in terms of cost and time reduction. Medical experts have a "reading fee" that they are charged to examine and diagnose any potential abnormality appearing on the radiology image. For example, a radiologist in California charges around $100 to $500 for the interpretation of an image. Radiology images are captured by a radiology machine technician who charged a heavy image reading fee. So, if we can build such a system that just takes the radiology image and answers the questions. Through such a VQA system, images will not be examined just for one medical problem but can be examined against all the usual questions related to a particular part of the patient's body. In more simple words, the radiology image will be examined against each of the given questions, and now the patient can be examined for all the medical problems that can be determined from the radiology image of that particular body part.

## 6.3 Purpose of Research

Such a system can also be beneficial for hospitals because medical imaging is the basis of modern-day medical diagnosis, and hospitals have specialist doctors that see the medical image and determine the disease. Specialist doctors are the major cost for the hospital, and many hospitals cannot even afford specialist doctors for each domain. For example, a radiologist usually takes 20, 000 dollars per month in the USA, and this amount is very huge for developing countries. Such VQA systems have a huge economic impact on patients as their costs are reduced due to the cost-cutting of hospitals. Such a system will

be a healthcare facility for the unprivileged. Medical VQA can be very beneficial both for the medical fraternity as well as for the patient in terms of cost and time reduction. Medical experts have a "reading fee" that they are charged to examine and diagnose any potential abnormality appearing on the radiology image. For example, a radiologist in California charges around $100 to $500 for the interpretation of an image. Radiology images are captured by a radiology machine technician and a heavy image analyst. analyst. reading fee. So, if we can build such a system that just takes the radiology image and answers the questions. Through such a VQA system, images will not be examined just for one medical problem but can be examined against all the usual questions related to a particular part of the patient's body. In more simple words, the radiology image will be examined against each of the given questions, and now the patient can be examined for all the medical problems that can be determined from the radiology image of that particular body part. So, this system will not only save 100 to 500 dollars for the patient but also examine the patient's radiographic image for each potential medical problem.

**6.4 Significance of the Study**

Such a system can also be beneficial for hospitals because medical imaging is the basis of modern-day medical diagnosis, and hospitals have specialist doctors that see the medical image and determine the disease. Specialist doctors are the major cost for the hospital and many hospitals even cannot afford specialist doctors for medical diagnosis. For this purpose, A medical VQA system has been developed for medical diagnosis.

## 6.5 Research Objectives

A system of this nature can also prove advantageous for hospitals, as medical imaging plays a crucial role in contemporary medical diagnosis, and hospitals often have specialists who examine medical images and diagnose conditions. Specialized physicians are a significant expense for hospitals, and some hospitals may not have the financial resources to hire specialists for every field.

In particular, the study has the following sub-objectives:

- Open-ended questions are those that allow for a wide range of possible answers and can be used to generate natural language responses for medical questions answers. Closed-ended questions, on the other hand, have a limited set of possible answers and can be used for tasks such as medical image classification, where the goal is to provide a specific answer or perform a specific action. These questions are useful for tasks where the system needs to identify specific information or perform a specific medical task.

- Our VQA system incorporates a multimodal input interaction, which means that it can process and understand multiple types of input, such as text and visual data, at the same time. This allows the system to better understand the question and the context of the image, and to provide a more accurate and relevant answer.

# CHAPTER 2.

# REVIEW OF LITERATURE

## 2.1 Preliminary Literature Review Objectives

General Visual Question Answering (VQA) is a multifaceted Artificial Intelligence (AI) research problem. This multidisciplinary field combines the domains of Computer Vision, Natural Language Processing, and Knowledge Representation & Reasoning. VQA datasets are composed of real-world images and some open or closed-ended questions against them. So, in VQA, primarily the merger of CV and NLP is used to understand the knowledge representation between both domains by transferring and correlating the knowledge space. VQA aims to answer the queries by extracting and then correlating the relevant visual information.

General VQA contains real-world scenarios with natural language questions, but in Medical Question Answering (MQA), we are provided with medical images of the internal structures of patients. Medical Imaging (MI) is the domain of producing images of the internal vision of the patient's body; it can benefit clinicians in disease diagnosis. The MI helps the doctor visualize and analyze the patient's physical conditions to diagnose and assess the irregularities. Different medical practitioners may interpret the same medical image differently.

Deep Learning Models (Goyal et al., 2017) are used to process complex data like medical images. Deep learning methods are available to extract relevant information from visual or textual data. In the context of medical imaging, a deep learning-based medical VQA should have two parts: the first part that can process medical images and extract relevant information, and the second part that can understand the question and answer by extracting the relevant information from the image. Such a medical-VQA system can help doctors in diagnosis, while patients can benefit from the medical-VQA system by getting a preliminary understanding of their body condition, cost-effective diagnosis, and can guide in choosing a more precise treatment plan.

When it comes to medical VQA, several aspects should be kept in mind. There are multiple types of images available in the medical domain. Different types of modalities exist for medical images like X-ray, Radiology, CT, MRI, etc. An image may contain two or more types of irregularities if it is diagnosed by different field experts. The unavailability of enough data samples for training deep learning models also poses a problem. Moreover, general VQA networks are fine-tuned on the medical images, but medical images possess a great amount of complexity as compared to daily life images that demand a specific VQA architecture. In addition to the dataset and VQA network problems, there is also a massive vocabulary difference between general and medical question answering.

Existing medical-VQA systems (Zhan et al., 2020; Abacha et al., 2018) are based on transfer learning and exploit the existing deep learning-based VQA system by fine-tuning the pre-trained models on the comparatively smaller medical datasets that are pre-trained-on the general datasets. However, there is a very great gap between the general and medical domains, both in terms of image structuring as well as the language of question-answer pairs (Litjens et al., 2017). We have reviewed the state-of-the-art medical VQA systems and identified the challenges faced by them.

## 6.6 Medical-VQA

In general, VQA's state-of-the-art deep learning model has been used but (Goyal et al., 2017) duly pointed out that more models exploit the language priors and ignore the visual information in images. VQA models initially employ classification, object detection, and segmentation to extract the relevant information, and knowledge-based reasoning is used to answer the questions, but in medical VQA due to the high complexity of medical data and relatively complex questioning. Med VQA is also risky because a patient's medication or treatment depends on the results that demand highly accurate results. VGGNet, ResNet, and Inception are used to extract visual information from images, while question-answer pairs related to the image are processed via Recurrent Neural Networks (RNN) such as Long Short-Term Memory (LSTM) (Greff et al, 2016) and Gated Recurrent Units (GRU) (Chung et al., 2014) are employed to encode the text representation and output the answer. However, the answers produced by these models are not complete sentences. Therefore, pre-trained GPT-3 (Brown et al., 2020) and BERT (Devlin et al., 2018) are used to

construct proper sentences. Fusion of the self-attention mechanism and another feature extractor may be used due to the versatility of this domain.

The parallel structures ResNet152 (K. He & Sun, n.d.; Kathiravan Srinivasan, Lalit Garg, Debajit Datta, Abdulellah A. Alaboudi, N. Z. Jhanjhi, Rishav Agarwal, 2021) and Gate Recurrent Unit (GRU) (Merri, 2013) were used to extract local features of the image. Its goal was to save spatial features from images in various dimensions. The basic three-channel images were then converted to separate-channel grayscale images and fed into the stacked GRU network to maintain the images' sequence feature information. Finally, the features extracted out of each layer of ResNet152 as well as the output of the GRU network were convolved to form complete image features. Moreover, the accuracy of predicted answers was not ideal. Although the model achieved cutting-edge performance and the results were less impressive.

(S. Liu et al., 2022) proposed an innovative bi-branched method for medical visual questionaries (BPI-MVQA) based on the parallel network and image retrieval that can be used in various classification methods for multiple types of training data for the VQA-Med. The first branch uses a transformer model, classified in parallel (Furfari(tony), 2002) to extract Image features. The second branch employs a method that retrieves image similarity and outputs the labels of relevant images as text descriptions. The pre-trained VGG16 network (Simonyan & Zisserman, 2015) was used in a novel method that removes the fully connected layer to output the feature representation of the image and then selects

the answer labels of relevant images by calculating the covariance of feature matrices of two images. This method enhances the precision of a portion of the test set data.

(Pan et al., 2021) proposed the MuVAM model, which effectively solves medical VQA tasks. It was divided into three modules. In the first module for feature extraction, two feature extraction methods were utilized to obtain input images as well as question representation. Second, this study suggested a multi-view attention module to maximize the implementation of semantic information. which included word-to-text (W2T) attention and image-to-question (I2Q) attention, which investigated the potential influence of the image and word on the question. Third, a compound loss module was proposed to train the model to improve MuVAM's accuracy. This research consisted of image-question-complementary (IQC) loss and classification loss. It was worth noting that the IQC loss used image representation and text semantics to collectively direct the question of the significance of learning to enhance the role of resemblance and weakened the difference in visual-text cross-modal attributes. This study's experiment corrected and completed the VQA-RAD dataset and created an enhanced dataset called VQA-RADPh to improve data quality Their results on these two datasets showed that MuVAM outperforms the state-of-the-art methods.

## 6.7 Datasets

Domain-specific datasets for the VQA system including image and related questions-answer pairs are the core for utilizing the deep learning-based computer vision and NLP

models. To the greatest of our knowledge, the following medical VQA datasets are currently available for public use: VQA-RAD (Lau et al., 2018), VQA-MED-2018 (Hasan et al., 2018), VQA-MED-2019 (Abacha, Datla, et al., n.d.), PathVQA (He et al., 2020), RadVisDial (Kovaleva et al., 2020), VQA-MED-2020 (Abacha, Datla, et al., n.d.), VQA-MED-2021 (Abacha, Datla, et al., n.d.; Abacha, Datla, et al., n.d.; Pelka et al., 2021), and SLAKE (Liu et al., n.d.). The following sections give an overview of the QA sets collection. Table. 3 shows the samples of medical images and QP pairs of various datasets, while Table 4. provides the summary of different datasets

*Table 4: Samples of medical images and QP pairs of various datasets.*

| VQA-Med-2018 |  | Q: how was this image taken?<br>A: Xr - plain film |  | Q: what imaging modality was used to take this image?<br>A:mr-t1w w/gadolinium |
|---|---|---|---|---|
| VQA-Med-2019 |  | Q: what abnormality is seen in the image?<br>A: inguinal hernia involving bl. |  | Q: what organ systems can be evaluated with this MRI? |

| | | | A: skull and contents. |
|---|---|---|---|
| PathVQA |  | Q: What is the end of the long bone expanded in?<br><br>A: Region of epiphysis |  | Q: What shows a large and tan mass while the rest of the kidney has a reniform contour?<br><br>A: Upper pole of the kidney |
| RadVisDial |  | Q: Fracture?<br><br>A: Not in the report. |  | Q: Pneumonia?<br><br>A: Yes |

| VQA-RAD |  | Q: What is the organ system? <br><br> A: Gastrointestinal |  | Q: What is abnormal in the gastrointestinal image? <br><br> A: Gastric Volvulus. |
|---|---|---|---|---|
| SLAKE |  | Q: What is the function of the rightmost organ in this picture? <br><br> A: Breathe | | |
| VQA-Med-2020 |  | Q: What is most alarming about this MRI? <br><br> A: focal nodular hyperplasia |  | Q: What abnormality is seen in the image? <br><br> A: Enhancing lesion right parietal lobe with surrounding edema |

### 2.3.1. VQA-Med-2018

VQA-Med-2018 (Hasan et al., 2018) is the first available public data source in the medical domain and was proposed in Image CLEF 20183. A semi-automatic method was used to generate the QA pairs from the captions. Initially, a rule-based question generation (QG) method generated potential QA pairs by simplifying sentences, identifying answer phrases, generating questions, and ranking candidate questions. The produced QA pairs were then manually checked by two trained human annotators in two passes. One pass guarantees semantic correctness, while the other makes sure of clinical relevance to related medical images.

### 2.3.2. VQA-RAD

VQA-RAD (Lau et al., 2018) is a radiology dataset that was introduced in 2018. The image dataset is balanced, with MedPix5 samples from the head, abdomen, and chest. The images proposed to healthcare professionals to collect unsupervised questions to start investigating the question in a realistic incident. Health professionals must create questions including both template and free-form structures. Following that, the QA pairs are manually validated and classified to analyze the diagnostic focus. There are two types of answers: open-ended and closed-ended. Despite its small size, the VQA-RAD data contains critical information regarding what a medical VQA process should be capable of answering as just an AI radiologist.

### 2.3.3. VQA-Med-2019

VQA-Med-2019 (Abacha, Hasan, et al., n.d.)  is VQAMed's second edition, released during the ImageCLEF 2019 challenge. VQA-Med-2019 was motivated by VQA-RAD (Lau et al., 2018) and addresses the four most common question groups: modality, organ system, plane, and abnormality. The questions in each group follow the trends of hundreds of normally asked and validated questions within VQA-RAD (Lau et al., 2018). The first three groups (modality, organ system, and plane) seem to be classification problems, whereas the fourth (abnormality) is an answer to the generation challenge.

### 2.3.4. RadVisDial

RadVisDial (Kovaleva et al., 2020) is an open-source dataset for visual discussion in radiology. The visual discussion, which includes several QA pairs, is thought to be a practical and challenging problem for a radiology AI system such as VQA. The images were chosen from MIMIC-CXR (Johnson et al., 2019) and offered a well-structured applicable report with annotations for fourteen labels for each image. The RadVisDial is composed of two datasets: gold-standard and silver-standard. The dialogues in the silver-standard group are generated synthetically from the plain text reports connected with every image. Each dialogue includes five questions chosen at random from a pool of thirteen possible questions.

Each dialogue includes five questions chosen at random from a pool of thirteen possible questions. The respective answer is extracted automatically from the data source and is

restricted to four options (maybe, not mentioned, yes, no). To maintain consistency, the gold-standard cluster collects dialogues from two experienced radiologists' discussions using detailed annotation standards. Only 100 images at random are labeled as the gold standard. The RadVisDial data source investigated a real-world AI scene task in the medical field. Furthermore, the group compared synthetical dialogue to real-world dialogue and did experiments to demonstrate the significance of contextual information. The patient's medical history was also introduced, which improved accuracy.

### 2.3.5.   PathVQA

PathVQA (He et al., 2020) is a data source that investigates VQA in pathology. The images with annotations are retrieved from digital resources (online libraries and electronic textbooks). The author created a semi-automated pipeline to generate the annotations into QA pairs, which are then manually examined and revised. What, when, where, how, whose, how many/how much, and yes/no are the seven categories of questions. Open-ended questions make up 50.2% of overall questions. The answers to the closed-ended "yes/no" questions are balanced, with 8,145 "yes" as well as 8,189 "no".

The questions are based on the American Board of Pathology's pathologist certification examination (ABP). As a result, it is an exam to validate the "AI Pathologist" through decision support. The PathVQA (He et al., 2020) dataset shows how medical VQA can be used in a variety of scenes.

### 2.3.6. SLAKE

SLAKE (Liu et al., n.d.) is a large dataset that includes both a structural medical and semantic labels knowledge base. The images are drawn from three open-source data sources (Kavur et al., n.d.; X. Wang et al., n.d.; Kavur et al., n.d.) and annotated by doctors. For visual objects, semantic labels for images focus on providing bounding boxes (detection) and masks (segmentation). The medical expertise is presented as a knowledge graph. and manually reviewed after being retrieved from Own Think. They are presented in the pattern of triplets like (Heart, Function, and Enhance blood flow).

There are 2,629 triplets in Chinese and 2,603 triplets in English in the dataset. The use of a knowledge graph enables exterior knowledge-based questions including disease prevention and organ function to be answered. The questions are gathered from medical experts and doctors by choosing pre-defined questions. The questions are then classified by type and balanced to prevent bias.

### 2.3.7. VQA-Med-2020

VQA-Med-2020 (Abacha, Datla, et al., n.d.) is the third edition of the VQA Med and was released as part of the ImageCLEF 2020 challenge. The images were chosen with the restriction that the prognosis was made based on the image information. The inquiries are as follows: Actually, addressing abnormality. A collection of 330 abnormality problems is chosen, with each problem making an appearance at a minimum of ten times in the dataset. Patterns are used to generate the QA pairs. Visual question generation (VQG) is presented to the medical domain for the first time in VQA-Med-2020. The task of the

VQG is to create natural language questions about the image content. The medical VQG dataset contains 1,001 radiology images with 2,400 questions. The dataset questions are created using a rule-based approach and manually revised based on the image descriptions.

### 2.3.8. VQA-Med-2021

VQA-Med-2021 (Abacha, Datla, et al., n.d.) has been accepted into the ImageCLEF 2021 problem. The VQA-Med-2021 is based on the same principles and training dataset as the VQA-Med-2020. The validation set and test set are completely new and have been manually examined by medical doctors.

### 2.3.9. Summary of VQA datasets

*Table 5:Summary of VQA datasets*

| Datasets | # of Images | # of QA Pairs | Source | QA Creation | Category Question |
|---|---|---|---|---|---|
| VQA 2.0 (Parikh§., 2017) | 204k | 614k | COCO (T. Y. Lin et al., 2014) | Selected Manually | • Object<br>• Sport<br>• Color<br>• Count |

| VQA-Med-2018 (Hasan et al., 2018) | 2866 | 6413 | PubMed Central Repository | Synthetic | • Yes\No Questions<br>• Location<br>• Finding<br>• Other Questions |
|---|---|---|---|---|---|
| VQA-RAD (Cheng et al., 2022; Lau et al., 2018) | 315 | 3515 | MRIs<br>- Chest X-rays<br>- Abdominal axial CTs<br>MedPix Repository<br>- Head axial single-slice | Natural | • Plane<br>• Modality<br>• Organ System<br>• Abnormality<br>• Object/Condition Presence<br>• Positional reasoning<br>• Color<br>• Size<br>• Attribute Other<br>• Counting<br>• Other |

| VQA-Med-2019 (Abacha et al., n.d.) | 4200 | 15292 | MedPix database: - 10 organ systems, 16 Planes and Various in 36 modalities. | Synthetical | • Organ System • Plane • Modality • Abnormality |
|---|---|---|---|---|---|
| RadVisDial (Gold-Standard) | 100 | 500 | MIMIC-CXR (Johnson et al., 2019): -posterior-anterior (PA) view of chest X-ray. | Natural | • Abnormality |
| RadVisDial (Silver-Standard) | 91060 | 455300 | MIMIC-CXR (Johnson et al., 2019): -posterior-anterior (PA) view for chest X-ray. | Synthetical | • Abnormality |
| PathVQA (He et al., 2020) | 4998 | 32799 | Electronic pathology textbooks from | Synthetical | • Shape • Color |

| | | | PEIR Digital Library | | • Appearance Location<br>• Etc |
|---|---|---|---|---|---|
| SLAKE<br>(B. Liu et al., n.d.) | 642 | 14000 | NIH Chest X-ray (Wang et al., n.d.), Medical Segmentation Decathlon(Simpson et al., 2019), CHAOS(Kavur et al., n.d.):<br>- Neck CTs<br>- Abdomen CTs/MRIs<br>- Chest X-rays/CTs<br>- Head CTs/MRIs<br>- Pelvic cavity CTs | Natural | • Quality<br>• Organ<br>• Knowledge Graph<br>• Position<br>• Modality<br>• Abnormality<br>• Plane<br>• Color<br>• Shape<br>• Size |
| VQA-Med-2020 | 4k | 4k | MedPix Repository | Synthetical | • Abnormality |

| | | | | | |
|---|---|---|---|---|---|
| (Abacha, Datla, et al., n.d.) | | | | | |
| VQA-Med-2021` (Abacha, Datla, et al., n.d.) | 5k | 5k | MedPix Repository | Synthetical | • Abnormality |

## 2.4. Feature fusion techniques

Feature fusion is also significant for VQA tasks because of the combination of image and text to give the correct answer. Fusion techniques have evolved from hierarchical co-attention models to bilinear pooling (BLP). Bilinear Pooling (BLP) is developed for feature fusion from modalities predominantly developed for VQA models. A bilinear (outer product) has outperformed simple vector operations (concatenation & element-wise addition/multiplication) on VQA benchmarks. These techniques work better alongside attention mechanisms (Winterbottom et al., 2020).

Multimodal Compact Bilinear (MCB) pooling fused the visual features with the textual features and then used these features for question answering. (Fukui et al., 2016) used

MCB pooling twice to answer the question and worked on the Visual7W dataset and the VQA challenge. Multimodal Local Perception Bilinear Pooling is proposed by (Lao et al. (2018), a novel multimodal feature fusion approach that retains the second-order interactions between visual and textual features with limited learning parameters. MLPB utilizes local perception mechanisms, transforming two high-dimensional raw features into multiple low-dimensional part features. They have reduced the computational cost by sharing the learning parameters of each local bilinear pool.

The fusing of the textual and visual features proves to be beneficial in VQA, but it leads to high computational complexity and reduces its applicability practically. (Yu et al., 2017) explore feature fusion to mitigate this problem. They developed Multi-modal Factorized Bilinear (MFB) Pooling for an effective combination of these features with co-attention mechanisms.

## 2.5. Fine-Tuned VQA Model

(Zhan et al., 2020) proposed a conditional reasoning framework for medical VQA where the conditional reasoning module learned a different set of reasoning skills for open-ended and closed-ended questions. However, the reasoning for open and closed-ended questions was learned through separate training, and open-ended answers also lacked reasoning ability. (Abacha et al., 2018) used Stacked Attention Network (SAN) and Multimodal Compact Bilinear Pooling (MCB). SAN performed relatively better in medical VQA. (Pan, data and 2009; Brown et al., 2020) explained how existing systems used pre-trained

models on a larger dataset and then fine-tuned them on a comparatively smaller medical image dataset. But there is a gap between general and medical images as well as the question-answering pairs (Litjens et al., 2017). These approaches are not sufficient enough in the medical domain.

Gong et al. (2021) used hierarchical feature extraction to capture multi-scale features of medical images, and data augmentation to deal with data limitation with a curriculum learning paradigm being used with label smoothing and ensemble learning. Most of the fine-tuned medical VQA models use pre-trained feature extractors that are trained on general image datasets, but medical datasets have different content as compared to general datasets like ImageNet. Such VQA systems regard the task as classification only instead of mapping the corresponding information into image and QA. A slight shift in the data distribution of images can result in regarding the performance when using pre-trained feature extractors and weights from the general domain.

To learn bi-level image features for medical VQA from limited data, (Li et al., 2022) possessed an effective model BiRL. It used sentence-level rationale to extract fine-grained representations (Pan et al., 2021) from input questions and images, and token-level rationale to build a fine-grained multi-model vector that guided the model to adaptively filter the insignificant semantic representations of coupled questions and images. The model learned fine-grained semantic features from a limited scale of medical VQA data after combining the two rationale modules.

Furthermore, this study introduced the LDSM loss, to minimize generalization errors, bonded by smoothly changed label-margin bound in long-tailed label dispersion in medical VQA data. Empirical evaluation results on basic benchmark datasets revealed that the model performed better than state-of-the-art models, based on the standard PathVQA dataset and the VQA-Rad dataset, the suggested model achieved 0.5434 and 0.7605 in accuracy, as well as 0.5288 and 0.7741 on F1-score, outperformed the state-of-the-art benchmark models.

For visual questions answering from radiology images, the CNN model for visual feature extraction and the Bidirectional Long Short-Term Memory (BiLSTM) model for feature extraction for textual data are the better choices. (Y. I. Jinesh Melvin, 2022) suggested a method that solved image classification problems using CNN and text classification problems using BiLSTM. This system aided in the feature extraction from radiology images and provided users with suitable answers to their questions, which should be both objective and descriptive.

A visualization method with greater accuracy projected the answers as a benchmark that demonstrated the corresponding area with different colors and enabled them to mention the answers together in the visual method for the relevant questions. The benefit of this study was that the answers were traced back for more precision and potential treatments. The method of BiLSTM is to assign suitable weights for the radiological image based on the similarity measures of the question-and-answer pair.

## 2.6. Attention-based models.

Attention mechanisms including self-attention, co-attention, and multi-head attention are widely used in transformer models for various text and image analysis tasks, such as image captioning (Devlin et al., 2018) and natural language models such as GPT-3 and BERT (A. Ben Abacha et al., 2019). (Devlin et al., 2018) introduced a new language representation model, BERT. It is designed to pre-train deep bidirectional representations from the unlabeled text by joint conditioning on both left and right contexts in all layers. It is to be fine-tuned with minimal changes. (Alsentzer et al., 2019) have discussed currently available language models ELMo and BERT to solve many NLP problems in a general context, but in the clinical domain, no pre-trained model yet exists. Med-BERT has released BERT models for domain-specific problems, particularly in the clinical domain. Brown et al. (2020) worked on scaling up the language model, particularly GPT-3, to achieve strong performances on the NLP dataset for question-answering, translation, etc.

The Stacked Attention Network (SAM) (Yang et al., 2016) uses the question feature as a query to apply attention to rank the image regions that are relevant to the answer. Relevant answers are generated progressively with the incursive attention and multi-layer architecture that make the network capable of querying the image multiple times. Such a VQA mechanism cannot be effective where the answer cannot be formed progressively. This network lacks the joint learning of multi-model input that is required in the medical

51

VQA due to its complex nature and the accuracy required. MMBERT (Khare et al., 2021) uses multi-head attention and tries to learn semantic representations. QA images and QA text are created by applying the Vision-language model to the images having related captions for QA text, but this network lacks the joint learning of multi-model input. The Encoder-decoder network (Yan et al., 2019) used CNN and BERT to get the image and QA textual representations. Image and textual representations are fused by using a co-attention and fed to the decoder for answer generation. The Hierarchical co-attention VGQ model (Lu et al., 2016) reasoned about the image and questioned attention jointly where interpretable image and question regions were co-attended to generate the answer. Bilinear attention networks (Kim, Jun, and Zhang, 2018) have formed two separate attention flows for each modality and ignore the multimodal inputs. Bilinear attention used the visual and textual information smoothly to predict the answer.

State-of-the-art attention mechanism frameworks that perform well on various vision-textual tasks, including Modular Co-attention Networks and Transformer (Vaswani et al., 2017), are not explored much for the medical VQA. Such networks offer a nice way to learn the richer representation jointly for multi-modal and multi-channel inputs. Hence, the usage of attention mechanisms in medical VQA still has huge potential to design a robust medical VQA system.

Due to the lack of medical professionals worldwide, the large number of cases causes mental and physical fatigue, resulting in human errors during diagnosis. In such cases,

getting an additional opinion can help to boost the decision maker's confidence. As a result, having a dependable visual question answering (VQA) system and providing a second opinion on medical fields becomes critical. Moreover, often these VQA systems in use today are designed to solve real-world problems and aren't specifically designed to handle medical images.

(Sharma et al., n.d.) created MedFuseNet, an attention-based multimodal for VQA on medical images that take into account the associated challenges. Besides having broken the problem statement into simpler parts and predicting the answer, MedFuseNet aimed to maximize learning while minimizing complexity. MedFuseNet tackled answer prediction in two ways: generation and categorization. This study used a comprehensive set of both qualitative and quantitative analyses to assess the performance of MedFuseNet. Results of experiments demonstrated that MedFuseNet outperformed state-of-the-art VQA approaches, and visualization of captured attention demonstrated the predictability of the MedFuseNet model's expected results.

Computer-aided diagnosis has the potential to alleviate some of the present states of medical resource imbalance, and medical images are increasingly being used in medically assisted prognosis. The CGMVQA was proposed by (Ren et al., 2020) for answering questions based on medical images, except for previous work, the model was not limited to a single disease and can be applied to a variety of medical images and different organs.

In particular, the ImageCLEF 2019 VQA-Med data was used and split into 5 groups to simplify the complex problem, and a comprehensive model, such as answer generation capabilities and classification, was proposed. Because of the limited availability of data, this study used data augmentation on images and tokenization on texts. A pre-trained ResNet152 model was used to extract feature representations, and a global average pooling approach was applied to unify the dimensions of such features. The segment and position embedding layers were added together to deal with texts. To avoid the warm-up optimizer, text, tokens, and image features were integrated as input to the model, which used the pre-layer-normalization multi-head self-attention transformer. To ensure that the model can be executed on a single GPU, the parameters were reduced, and the embedding weight was shared.

Getting a Surgical-VQA system like a trusted "second opinion" could serve as a backup and relieve the burden on medical experts in answering these questions. The research of VQA for surgical processes has been limited due to the absence of annotated medical data and the involvement of domain-specific aspects. (Seenivasan et al., 2022) designed a Surgical-VQA that answers surgical procedure questionnaires based on the surgical scene. By introducing two Surgical-VQA data sources with classification and statement answers, which extend the MICCAI (Medical Image Computing and Computer Assisted Intervention) endoscopic vision challenging problem 2018 dataset. Surgical-VQ performed by vision-text transformer models and residual MLP-based VisualBert encoder used that imposes interaction among text and visual tokens, enhancing classification-

based answered performance. By incorporating a cross-token sub-module, the VisualBert encoder surpasses the vision-text attention encoder model. The impact of the model's performance on the input image patches and the integration of temporal visual features are also discussed. During Surgical-VQA answers to the less-complex questionnaire; from an application perspective, this study allowed for the possibility of integrating open-ended questions in which the model was trained to respond to surgery-specific complex questions. In addition, study the impact of the input image patches and temporal visual attributes on model performance in classification and sentence-based answering.

To train the VQA system, only images and questions were used, and they were categorized based on the results of the "[CLS]" position. Except for abnormality, the classification model applies to other categories. Data imbalance has an impact on this model. Answers to an abnormality category were generated using questions, images, and masked answers. By looping, the generative model predicted the sequence. The number of data limits the generative model's output. Strict accuracy, semantic similarity, and inward matching were adopted as evaluation metrics. On ImageCLEF 2019 VQA-Med data, the proposed model produced state-of-the-art results with a 0.659 BLEU score, 0.640 accuracies, and a 0.678 WBSS score.

## 2.7. Summary of the proposed VQA systems

*Table 6 provides a summary of the proposed VQA systems.*

| Model | Language Encoder | Image Encode | Fusion | Attention | Output Mode | Language score | Acc |
|---|---|---|---|---|---|---|---|
| **VQA-MED-2018** | | | | | | | |
| FSST (Imane Allaouzi et al., n.d.) | Bi-LSTM | VGG16 | Concatenation | N0 | Classification | 0.054 | |
| TU (Zhou, Kang, Notes), & 2018, n.d.) | Bi-LSTM | Inception-Resnet-v2 | Attention mechanism | Yes | Classification | 0.135 | |
| UMMS (Peng et al., n.d.) | LSTM | ResNet-152 | MFB using Co-attention | Yes | Classification | 0.162 | |
| JUST (Talafha et al., n.d.) | LSTM | VGGNet | Concatenation | No | Generation (LSTM) | 0.016 | |
| Chakri (Ambati et al., n.d.) | GRU | VGG16 | Element-wise multiplication | No | Generation (GRU) | 0.188 | |
| NLM (Rajaraman et al., 2018) | LSTM | VGG16 | SAN | Yes | Classification | 0.121 | |

| VQA-MED-2019 | | | | | | | |
|---|---|---|---|---|---|---|---|
| MedFuseNet (Sharma et al., n.d.) | LSTM | ResNet-152 | MedFuseNet | Yes | Classification & Generation (LSTM) | 0.27 | 0.789 |
| UMMS (Shi et al., n.d.) | Bi-LSTM | ResNet-152 | MFH with Co-attention | Yes | Classification | 0.593 | 0.566 |
| Techno (Z. Lin et al., n.d.) | LSTM | VGG16 | Concatenation | No | Classification | 0.486 | 0.462 |
| IBM Research AI (Kornuta et al., n.d.) | LSTM | VGG16 | Attention-Mechanism Classification | Yes | Question Classifier | 0.582 | 0.558 |
| MMBERT (Kornuta et al., n.d.) | Transformer | ResNet-152 | Multi-Head-Attention | Yes | Generation (Transformer) | 0.69 | 0.672 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Hanlin (Al-Sadi et al., n.d.) | Transformer | VGG16 | MFB with Co-attention | Yes | Classification | 0.644 | 0.62 |
| JUST19 (Al-Sadi et al., n.d.) | LSTM | VGG16 | None | No | Classification/Generation (LSTM) | 0.591 | 0.534 |
| KEML (Zheng et al., 2020) | Transformer | VGG16 | Block | No | Classification | 0.912 | 0.938 |
| TUAI (Zhou, Kang, Notes), & 2019, n.d.) | Transformer | Inception-Resnet-v2 | QC-MLB | Yes | Classification | | 0.603 |
| LIST (I Allaouzi et al., n.d.) | LSTM | DenseNet-121 | Concatenation | No | Generation (LSTM) | 0.583 | 0.556 |
| **VQA-MED-2020** | | | | | | | |
| NLM (Notes) & 2020, n.d.) | None | ResNet-50 | None | No | Classification | 0.441 | 0.4 |

| kdevqa (Umada et al., n.d.) | Transformer | VGG16 | GLU | No | Classification | 0.35 | 0.314 |
|---|---|---|---|---|---|---|---|
| Shengyan (S. Liu et al., n.d.) | GRU | VGG16 | None | No | Generation (GRU) | 0.412 | 0.376 |
| HCP-MIC (Chen et al., n.d.) | Transformer | BBN-ResNeSt-50 | None | No | Classification | 0.462 | 0.426 |
| Bumjun_jung (Jung et al., n.d.) | Transformer | VGG16 | MFH with Co-attention | Yes | Classification | 0.502 | 0.466 |
| **VQA-MED-2021** | | | | | | | |
| TAM (Y. Li et al., 2021) | LSTM | Modified ResNet-34 | MFB with co-attention | Yes | Classification | 0.255 | 0.222 |
| Lijie (J. Li et al., 2021) | Transformer | VGG8 | MFB with co-attention | Yes | Classification | 0.352 | 0.316 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| SYSU-HCP (Gong et al., 2021) | None | ResNets, plus HAGAP, and VGGNet | None | No | Classification | 0.416 | 0.382 |
| Sheerin (Sitara et al., n.d.) | LSTM | VGGNet | Element-wise-Multiplication | No | Generation (LSTM) | 0.227 | 0.196 |
| IALab_PUC (Schilling et al., n.d.) | None | DenseNet-121 | None | No | Classification | 0.276 | 0.236 |
| **VQA-RAD** | | | | | | | |
| MMQ (Do et al., 2021) | LSTM | MMQ | BAN/SAN | Yes | Classification | | 0.67 |
| HQS (Gupta et al., n.d.) | Bi-LSTM | InceptionResnet-v2 | Concatenation | No | Classification | 0.411 | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| CPRD (B. Liu et al., 2021) | LSTM | ResNet-8 | BAN | Yes | Classification | | 0.727 |
| QCR (Zhan et al., 2020) | LSTM | MEVF | BAN/SAN | Yes | Classification | | 0.716 |
| HQS (Gupta et al., n.d.) | Bi-LSTM | Inception-Resnet-v2 | Concatenation | No | Classification | 0.411 | |
| MEVF (Nguyen et al., 2019) | LSTM | MEVF | BAN/SAN | Yes | Classification | | 0.627 |
| **SLAKE** | | | | | | | |
| CPRD (B. Liu et al., 2021) | LSTM | ResNet-8 | BAN | Yes | Classification | | 0.821 |
| **PathVQA** | | | | | | | |

| MMQ (Do et al., 2021) | LSTM | MMQ | BAN/SAN | Yes | Classification | | 0.488 |
|---|---|---|---|---|---|---|---|
| MedFuseNet (Sharma et al., n.d.) | LSTM | ResNet-152 | MedFuseNet | Yes | Classification/ Generation (LSTM) | 0.605 | 0.636 |

## 2.8. Summary

Most of the existing work is focused on general VQA and medical VQA is an emerging field so the system should be designed to particularly address the problems in these medical images. Transfer learning, or fine, is used extensively in deep learning. However, A change in data distribution can affect the performance when pre-trained weights. and feature extractors are trained on general image datasets like ImageNet and MS COCO. In addition to this, there is a greater co-occurrence diversification of words in both medical text and general domain text. These factors motivate the need for an attention-based model and a time-based transformer to learn effective representations. From the above details, we can conclude that medical VQA has huge potential in the medical domain. However, only transfer learning and fine-tuning of the deep learning models will not be enough given the complexity and life-critical nature of the domain.

# CHAPTER 3.

# METHODOLOGY

## 3.1 Overview of the Research Problem

Medical VQA can be very beneficial both for the medical community as well as for the patient in terms of cost and time reduction. Medical experts have a "reading fee" that they are charged to examine and diagnose any potential abnormality appearing on the radiology image. For example, a radiologist in California charges around $100 to $500 for the interpretation of an image. Radiology images are captured by a radiology machine technician who charged a heavy image reading fee. So, if we can build such a system that just takes the radiology image and answers the questions. Through such a VQA system, images will not be examined just for one medical problem but can be examined against all the usual questions related to a particular part of the patient's body. In more simple words, the radiology image will be examined against each of the given questions, and now the patient can be examined for all the medical problems that can be determined from the radiology image of that particular body part. So, this system will not only save 100 to 500 dollars for the patient but also examine the patient's radiographic image for each potential medical problem.

Such a system can also be beneficial for hospitals because medical imaging is the basis of modern-day medical diagnosis, and hospitals have specialist doctors that see the medical image and determine the disease. Specialist doctors are the major cost for the hospital, and many hospitals cannot even afford specialist doctors for each domain. For example, a radiologist usually takes 20, 000 dollars per month in the USA, and this amount is very huge for developing countries. Such VQA systems have a huge economic impact on patients as their costs are reduced due to the cost-cutting of hospitals. Such a system will be a healthcare facility for the unprivileged. Medical VQA can be very beneficial both for the medical fraternity as well as for the patient in terms of cost and time reduction. Medical experts have a "reading fee" that they are charged to examine and diagnose any potential abnormality appearing on the radiology image.  For example, a radiologist in California charges around $100 to $500 for the interpretation of an image. Radiology images are captured by a radiology machine technician and a heavy image analyst. analyst. reading fee. So, if we can build such a system that just takes the radiology image and answers the questions. Through such a VQA system, images will not be examined just for one medical problem but can be examined against all the usual questions related to a particular part of the patient's body. In more simple words, the radiology image will be examined against each of the given questions, and now the patient can be examined for all the medical problems that can be determined from the radiology image of that particular body part. So, this system will not only save 100 to 500 dollars for the patient but also examine the patient's radiographic image for each potential medical problem.

## 3.2 Methodology

### 3.2.1. Multimodal Interaction through Transformers

As shown in Figure 1, this study proposes a multimodal Transformer VQA model which is a composition of three modality forms including medical images and question-answer pairs input. EfficientNet is used as the feature extract to get the image embedding, while Bidirectional Encoder Representations from Transformers (BERT) is used to extract the text embedding from the question and answer. Multimodal representation interaction is performed with two transformers having multiple layers where one transformer is for inter-modality input representation and the other Transformer is for intra-modality input representations. Trilinear Interaction Attention (TrI-Att) fusion Transformer for the Inter-modality representations interaction, while self-attention for the intra-modality representation interaction (Yu et al., 2019). Then a guided-attention is used to learn inter-association between the question hyperedge and knowledge (image and answer modality) hyperedge. Finally, representations obtained from the attention layer are fed into a feed-forward layer (FFL) to make a joint representation and that representation is fed to the answer predictor for producing the answer. Guided attention and self-attention blocks are used within the transformer encoder for the prediction of the answer with the answer predictor in an output layer. Figure 2 shows the inference stage of the proposed VQA network.

## 3.3 Embeddings Extraction of Single Modality

### 3.3.1. Images Embeddings

The image embedding is obtained from the regional visual features extractor called EfficientNet (Koonce, 2021). In respect of specifications, it generates the vector of the $d_v$ dimensions for every item. The embedding matrix $V \in R^{vxd_v}$ is used to represent an image



*2: Research methodology of Training VQA-System: Multimodal Transformer VQA model, the composition of three modality forms including medical images and question-answer pairs input. EfficientNet is used for feature extraction*

66

*to get the image embedding, while Bidirectional Encoder Representations from Transformers (BERT) is used to extract text embedding from question and answer. Trilinear Interaction Attention (TrI-Att) fusion Transformer for the Inter-modality representations interaction, while self-attention for the intra-modality representation interaction. Then a guided-attention is used to learn inter-association between the question hyperedge and knowledge (image and answer) hyperedge. Representations obtained from the attention layer are fed into FFL to make joint representation and that representation is fed to the answer predictor for producing the answer. Guided attention and self-attention blocks are used within the transformer encoder for the prediction of the answer with the answer predictor in an output layer.*

of v objects. The EfficientNet is the convolution neural network design and scalability methodology that utilizes the compounded coefficient to adjust the depth, width, and resolution dimensions evenly. The Efficient Net scalability technique consistently increases network performance. breadth, depth, and resolution with a set of preset scalability coefficients contrasting with the standard practice of adjusting various parameters randomly (Tan et al., 2019). Figure 3 shows the layered architecture of EfficientNet.



*Figure 3: EfficientNet Architecture*

### 3.3.2 For embeddings of questions and answers

The absence of sufficient training data is one of the main problems facing NLP. Although there is a vast quantity of text data available overall, we must divide it up into many different fields to build task-specific datasets. We only have a few thousand to a few hundred thousand human-labeled training examples after doing this. Unfortunately, deep learning-based NLP models need a lot more data to function well; they significantly improve when trained on millions or billions of annotated training instances. Researchers

67

have created a variety of methods for training general-purpose language representation models using massive amounts of unannotated content on the web to fill up this data gap (this is known as pre-training). When working with issues like the question.



*Figure 4: Inference stage for testing VQA-System: In this section, the composition of two modality forms including medical images and question input. EfficientNet is used for feature extraction to get the image embedding, while Bidirectional Encoder Representations from Transformers (BERT) is used to extract text embedding from the question. Trilinear Interaction Attention (TrI-Att) fusion Transformer for the Inter-modality representations interaction, while self-attention for the intra-modality representation interaction. Representations obtained from the attention layer are*

answering and sentiment analysis, for example, these general-purpose pre-trained models can subsequently be fine-tuned on smaller task-specific datasets. When compared to starting from zero and training on smaller task-specific datasets, this method significantly increases accuracy. In the deep learning field, BERT, a relatively new addition to these techniques for NLP pre-training, generated a stir because it demonstrated cutting-edge outcomes in a range of NLP tasks, including question answering.

We use the BERT (Bidirectional Encoder Representations from Transformers), which is a state-of-the-art model used to deal with multiple NLP tasks. It was developed by Google (Devlin et al., 2018) to fine-tune our textual extraction. BERT is dependent on a Transformer (the attention mechanism that learns contextual relationships between words in a text). An encoder to read the text input and a decoder to create a prediction for the task make up a basic Transformer. Since the objective of BERT is to produce a language representation model, just the encoder portion is required. A series of tokens that are first transformed into vectors and then processed by the neural network make up the input to the BERT encoder. Specifically, the content is initially converted into Word Piece embedding. Each embedding is then projected onto $R^{d_q}$ or $R^{d_a}$, for question and answer, respectively, by fine-tuning. The final representation of the question with a maximum length of q is $Q \in R^{q \times d_q}$, and the equivalent for the answers is $A \in R^{a \times d_a}$.

Word embeddings are the representation of words in a machine-understandable numeric format. The simplest illustration is (Yes, No) expressed as (1, 0). But this may not be the most effective technique to represent words and sentences when working with vast texts and corpora. The co-occurrences of terms and their probabilities are crucial for huge corpora. Word Piece Embeddings of 30,000 token vocabularies are used by BERT. Every sequence always starts with a particular classification token as the first token ([CLS]). For classification tasks, the last hidden state matching this token is used as the aggregate sequence representation. A single sequence has several sentence pairings. Sentences are distinguished by a unique token ([SEP]), and each token is given a learned embedding indicating whether it belongs to sentence A or sentence B. The input representation for a particular token is built by adding the relevant token, segment, and position embeddings. These inputs are used to pre-train BERT for tasks like next-sentence prediction and masked language modeling. The architecture of BERT, however, makes it inappropriate for unsupervised tasks like clustering and information retrieval via semantic search as well as for tasks like semantic similarity search. BERT has revolutionized several NLP applications.

## 3.4  Theoretical Procedure

### 3.4.1. Inter-modality Representations with Trilinear Interaction Attention

The visual questioning and answering attention method is a successful technique. And attention mapping for the trilinear inputs is calculated through PARALING deconstruction with respect to the trilinear feature fusion. Nevertheless, the outcome is

just a classification combined matrix. We proposed Trilinear Interactive Attention to improve every individual modality (image, question, and answer) depiction by merging other modalities (TrI-Att). Furthermore, while self-attention cannot fuse different modalities, it can improve the interaction between them (Yu et al., 2019).

We design trilinear interaction. Attention to projecting the single-modality embeddings into an inter-modality improved environment for improved cross-modality data fusion as shown in Figure 3. The attention mechanism was influenced by how the human brain manages attention. Think of attending a party. Even if your name is lost in the background noise, you can still hear it being said on the other side of the room. Your brain can filter out all unimportant information and concentrate only on the things it deems significant. Queries, keys, and values aid in facilitating attention in transformers. A key is a label for a word that is used to identify one word from another. Query, look through all the keys that are offered, and choose the one that matches the best. It so signifies a live request for a specific piece of information. Value, key, and values are always presented in pairs. When a query matches a key, the word's value rather than the key itself spread. The knowledge that a word conveys is its value.

Let S = V, Q, and A be the multimodal information collections from the previous section (single modality embedding extraction). To begin, we present the $M \epsilon R^{vxqxa}$ attention map, which is primarily generated using matrix multiplication as well as sum-based down sampling. The following is a detailed calculating procedure:

$$M = softmax(\frac{\sum dv \sum dq \sum da V \otimes Q \otimes A}{\sqrt{d}})$$

where d is the arithmetic mean of $d_v$, $d_q$, and $d_a$, and softmax is an operation that normalizes all of the elements in M. Second, the attention map f is combined with the fusion of preliminary single-modality as follows:

$$fv = \sum q \sum a \, MV = TrI - Att_v(V, Q, A) \ \ldots\ldots \text{ Eq. 1.}$$

Here, we use image representation V as an example (questions and answers are identical), and the fusion operation is comparable to Eq. 1. By adding a linear mapping with every single-modality representation, we further improve the robustness using the attention mechanism (Vaswani et al., 2017). In particular, the inter-modality fusion calculation is as follows:

$$f_V = |_i^{N_h}| \ TrI\text{-}Att_V{}^i \ (VW^i{}_v, QW^i{}_Q, AW^i{}_v) \ \ldots\ldots \text{ Eq. 2.}$$

where the multi-head linear mappings $W_V{}^i$, $W_Q{}^i$, and $W_A{}^i$ are shared by the three types of representations. The number of heads is Nh. The symbol $||$ denotes the concatenation of every multi-head. Likewise, the fusion representations of the questions and answers are as follows:

$$f_Q = |_i^{N_h}| TrI\text{-}Att_Q{}^i \ (VW^i{}_v, QW^i{}_Q, AW^i{}_v)$$

$$f_A = |_i^{N_h}| TrI\text{-}Att_A{}^i \ (VW^i{}_v, QW^i{}_Q, AW^i{}_v) \ \ldots\ldots \text{ Eq. 3.}$$

The next step is a feed-forward network with a residual connection that is completely connected.

### 3.4.2. Intra-modalities representation with Self-Attention

To represent an intra-modality relationship, we would use the Transformers encoder (Vaswani *et al.*, no date). The multi-head self-attention method is used, accompanied by the feed-forward network with the remaining link. We use the Transformer encoder (Vaswani et al., 2017) to record the intra-modality (relationships within the modalities). After deploying a feed-forward (FF) network with residual connection, we deploy a multi-head self-attention mechanism. With the input potential $X \in R^{n \times d}$, as a result of the multi-head self-attention works as,

$$|_i^{N_h}|Self\text{-}Att_M(X) = |_i^{N_h}| \; softmax \left(\frac{XX^T}{\sqrt{d}}\right) XW^i_M \; \ldots\ldots \; Eq. \; 4.$$

where the projection matrix for $W^i_M \in R^{n \times d_h}$ a particular modality M inside its head This method can increase the reliance on long distances among the multiple modes while weakening the negative influence on the outcome to some extent.

### 3.4.3. Self Attention

Consider how the sequence of tokens is fed into the attention-pooling so that the same sequence of tokens serves as inquiries, keys, and responses. Every query, in particular, responds to every key-value pair and produces a single attention output. Self-attention (Lin *et al.*, 2017), also known as intra-attention, is performed because the searches, keys, and values all originate from the same location. Figure. 4 shows the Connectionist scheme of self-attention.

The searches, keys, and values in the self-attention are all *(n X d)* matrices. Considering scaled dot-products attention, which multiplies the *(n x d)* matrix by the *(d x n)* matrix, then multiplies the resultant *(n x n)* matrix by the *(n x d)* matrix, as a consequence, the computational complexities of self-attention are $O(n^2d)$. Every token is immediately related to any other token via self-attention. As a result, parallel processing is possible with O (1) sequential operations, and the maximum path length is also O (1).

Self-attention abandons sequential processes in favor of parallel computations by processing tokens of the sequence one by one on a regular basis. By introducing positioning encoding to input representations, we would inject absolutely or relatively positional information to leverage sequence order information. It is possible to learn and fix positional encodings. We would describe the fixed positional encodings based on the sine and cos functions in the following (Vaswani et al., n.d.).



Self attention

*Figure 5: Self-Attention*

Self-Attention enhances the semantics of the intermediate conceptions by allowing attention to the intra-modality and inter-modality aspects. It entails translating query vectors to weighted combinations of value-vectors, with the weights produced by multiplying the query and key vectors' dot product. In matrices Q, K, and V, query, key, and value vectors are all expressed simultaneously. The inversion of the dot-product of Q is $d_k$, wherein $d_k$ is a dimension of the question and key vectors. Several self-attentions (multi-head attention) are performed simultaneously by an encoder, and the output is concatenated. By responding to multiple representation subspaces at various positions, multi-head attention improves representation.

During this phase, the input image will be fed to Efficient Net, while questions and answers will be gone through pre-trained BERT to generate image and text embedding respectively. Then, embedding will be passed through trilinear attention to combining multimodal input. After that, individual and combined input will be passed through residual MLP followed by the self-attention layers. The output of the last residual MLP will be fed to the Transformer encoder and eventually to the output layer to predict the answer. During the interface, only images and questions will be fed as input, and answers will be output.

### 3.4.4. Feed-Forward Network

The feed-forward network (FFN) module implements a point-wise non-linear transformation on the extracted feature of the self-attention module using a two-layer

MLP model. The transformed features $F \in R^{n \times d}$ are extracted as follows from the input features $X \in R^{n \times d}$:

$$F = FFN(X) = ReLU\ (XW_i + b_1)W_2^T + b_2 \ \ldots\ldots \text{Eq. 5}.$$

Here the $W_1,\ W_2 \in R^{D \times 4D}$.

A feed-forward neural network (Figure. 5) is a type of artificial neural network in which there is no cycle in the connections between the nodes. A recurrent neural network, in which particular paths are cycled, is the reverse of a feed-forward neural network. Since the input is only processed in one direction, the feed-forward model is the simplest type of neural network. Although the data may flow via several buried nodes, it always proceeds forward and never backward. A single-layer perceptron is a common example of a feed-forward neural network in its most basic configuration. Several inputs are introduced into the layer in this model and multiplied by the weights. The weighted input values are then summed together to produce a total. The value produced is frequently 1, and if the sum of the values is below the threshold, the output value is -1. The threshold is typically set at zero. In classification tasks, the single-layer perceptron is a crucial feed-forward neural network model. Single-layer perceptron can also contain some features of artificial intelligence.

The neural network may compare the outputs of its nodes with the desired values using a property known as the delta rule, which enables the network to train its weights to create more accurate output values. This learning and training procedure results in gradient descent. Although the process of updating weights in multi-layered perceptron is almost

comparable, it is more formally known as back-propagation. In these circumstances, the network's hidden layers are each changed in accordance with the output values generated by the output layer.



*Figure 6: Samples feed-forward network of Multi-Layer-Perceptron (MLP)*

## 3.4.5. Transformer Encoder

Transformer is a revolutionary architecture that is introduced in the paper (Vaswani et al., n.d.) It employs the attention mechanism we already saw, as the title suggests. Transformer is an architecture for converting one sequence into another, much like LSTM. A novel guided and self-attention mechanism is used in this Transformer encoder network.

### 3.4.6. Guided Attention Blocks

We first embed an answer hyperedge and a question hyperedge in the manner described below to understand inter-association among two hypergraphs: $h^{[.]}$ is a $\epsilon^{[.]}$ hyperedge through and $e^k = \Phi_k^0 f_k (h_k) \in R^d$, while $e^q = \Phi_q^0 f_q (h_q) \in R^d$. Here, the linear projection function and the hyperedge embedding function are both present. We use a straightforward concatenation process of node representations inside a hyperedge as $f_{[.]}$ even though the design and implementation of $f_{[.]}$ are not restricted (e.g., any pooling process or any trainable neural networks). The representations of the same hypergraph's hyperedges (such as $e^k$ and $e^q$) are crammed into a matrix called $E^q$ and $E^k$.

As just a query and key-value pair, respectively, we describe the answer hyperedges $E^k$ and the question hyperedges $E^q$. All representation matrices $W_{[.]}$ $R$ $d^{dv}$ are trainable parameters, so we set a query $Q_k = E^K W_{Qk}$, a key $E_q = E^q W_{kq}$, and a value $V_q = E^q W_{Vq}$. When the query, key, and value have been used, the scaled dot product attention is calculated as Attention $(Q_k, K_q, V_q) = Vq$ softmax $(\frac{Q_k Kq_T}{\sqrt{d_v}})$ $Vq$, where $d_v$ is the dimension of the query and the key vector. Additionally, guided attention is decided to carry out correspondingly question hyperedges as the query and answer hyperedges as key-value pairs: $(Q_k, K_q, V_q)$.

### 3.4.7. Self-Attention Blocks

The only distinction between guided attention and self-attention is that during self-attention, the same input is used for both the key-value query. As an illustration, we fixed

the query, key, and value based on the information hyperedges $E_k$, and Attention is used

to perform the information hyperedges' self-attention $(Q_k, K_k, V_k)$. Self-attention is carried

out similarly for question hyperedges Eq: $(Q_q, K_q, V_q)$. The self-attention block and the

guided-attention block were constructed with the transformer's standard architecture, with

each block consisting of a single FF layer, layer normalization, and each attention

function. The representations of answer hyperedges and question hyperedges are modified

and finally concatenated to a separate vector representation as $z_k \in R^{d_v}$ and $z_q \in R^{d_v}$,

respectively, by transferring the self-attention blocks and guided-attention blocks in order.

### 3.4.8. Answer Predictor

To make a combined representation z to predict the answer, we initially combine the

representations zkand zq acquired from the attention blocks and feed them into a single

FF layer (for example, $R^{2d_v} \mapsto R^w$). Then, we take into account two different answer

predictor types: the MLP and the answer predictor based on similarity. For visual

question-and-answer challenges, the MLP, *or p = ψ(z)*, is frequently used.

For the answer based on similarity, we determine the dot product similarity $p = zC^T$ among

the input z and the answer candidate set $C \in R^{|A| \times w}$ w, here |A| is the total number of

candidate answers and w is the representational dimension for every answer. The answer

from answer candidates that most resemble the joint representation is chosen as the

answer. We don't annotate the ground-truth logic paths during training; instead, we only

use supervision from QA sets. Cross-entropy in between prediction p and the ground truth t is used as a loss function to achieve this.

## 3.5 Data Analysis

Firstly, we have PathVQA (He et al., 2020). This dataset has a total of 4,998 images and 32,799 pairs of question answers (QA pairs). All questions and answers are about color, location, appearance, and shape. The second source of datasets is the MedPix database, which contains head axial single-slice CTs, chest X-rays, and abdominal axial CTs. This dataset has 315 images with 3515 QA pairs. The nature of these question answers is natural, and it talks about the color size and position etc. (Abacha et al. 2019). The third source of datasets is the MedPix database, which contains 4,200 images and 15,292 QA pairs. The nature of this available dataset is synthetical, and it talks about the abnormality, plane, organ system, and modalities. (Abacha et al., 2019). The fourth source of datasets is the chest X-ray posterior-anterior (PA) view, which contains 91,160 images and 455,800 QA pairs. The nature of this available dataset is synthetical, and it talks about the abnormality. (Kovaleva et al., 2019) The final source of datasets is the MedPix database, which contains 5000 images and 5000 QA pairs. The nature of this available dataset is synthetical, and it talks about the abnormality. (Abacha et al., 2019).

## CHAPTER 4.

## RESULTS

### 4.1 Experiments

For the effective initialization of the proposed model, pre-training is carried out on the PathVQA and VQARAd datasets where the pre-training data format is a triple comprising the image, question, and answer for the trilinear model of the VQA system. While pre-train of Trilinear Interaction Transformer on the PathVQA and VQARAd, the question-and-answer pair is masked with a probability of 15%. In these masked tokens, sign [MASK] is substituted for 80% of them, 10% are preserved, and the remaining 10% are changed to random tokens. For the downstream task of VQA, the whole network is then trained on the med-VQA 2019 and Med-VQA 2022 datasets.

### 4.2 Hardware Set-up

In this study, we use PyTorch, which is an open-source machine learning library used to provide a dynamic computational graph, which allows users to change the graph on the fly during runtime, as opposed to a static computational graph. PyTorch also includes support for CUDA, which allows for the use of GPUs to accelerate computations to execute all of our experiments, so we are reliant on the API's available optimizations. We time the inference on three different hardware platforms, each of which corresponds to a

different use cases, to acquire more accurate timings:

One NVIDIA 1080 Ti GPU with a peak performance of 12 TFLOP/s. This training accelerator is standard.

A CPU with 64GB memory capacity and 9th generation Intel Core i7. This is a typical data center server that processes feature extraction on incoming image feeds. Using MKL and AVX2 instructions, PyTorch is well suited for this arrangement (16 vector registers of 256 bits each).

To simulate typical use cases on the GPU, we run timings on big image batches. We fine-tune to setting the parameters for the optimal performance on each GPU platform. We employ the largest strength batch size that resides in memory. To simulate a situation where multiple threads are processing various input image streams; we estimate inference time on CPU systems in a thread. Since it is challenging to separate the effects of the hardware and software, we test several network optimization techniques using common PyTorch tools (the just-in-time compiler, and various optimization profiles).

## 4.3 Datasets and Evaluation Metrics

### 4.3.1.  Dataset:

As we know in deep learning models, we need an ample amount of data to train our models and to get the accuracy and performance of the model. We have different datasets which are publicly available to address our problem. Firstly, we have PathVQA (He et al., 2020). This dataset has a total of 4,998 images and 32,799 pairs of question answers (QA pairs).

Electronic Pathology Textbooks PeIR Digital Library is the source of our dataset. We have synthetical questions and answers in nature. All questions and answers are about color, location, appearance, and shape. The second source of datasets is the MedPix database, which contains head axial single-slice CTs, chest X-rays, and abdominal axial CTs. This dataset has 315 images with 3515 QA pairs. The nature of these question answers is natural, and it talks about the color size and position etc. (Abacha et al. 2019). The third source of datasets is the MedPix database, which contains 4,200 images and 15,292 QA pairs. The nature of this available dataset is synthetical, and it talks about abnormality, plane, organ system, and modalities. (Abacha et al., 2019). The fourth source of datasets is the chest X-ray posterior-anterior (PA) view, which contains 91,160 images and 455,800 QA pairs. The nature of this available dataset is synthetical, and it talks about the abnormality. (Kovaleva et al., 2019) The final source of datasets is the MedPix database, which contains 5000 images and 5000 QA pairs. The nature of this available dataset is synthetical, and it talks about the abnormality. (Abacha et al., 2019). All the above-mentioned datasets are available publicly and are very useful for our model for training.

Then the whole is trained on VQA-Med-2019 and VQA-Med-2020 and the performance is evaluated on various evaluation metrics such as Bilingual evaluation understudy (BLEU), precision, recall, and accuracy. I employ a modified version of the accuracy metric from the broad domain VQA task that only takes into account exact matching between a question and answer. I determine the scores for each category of questions

along with the overall accuracy scores. The word overlap-based similarity between a model-generated answer and the ground truth response captured using BLEU (Janssens et al., n.d.), which makes up for the specificity of the accuracy metric. The BLEU metric's general methodology and available resources are roughly the same as the task from the previous year (Papineni et al., n.d.).

## 4.4 VQA-Med-2019

### 4.4.1. Question Patterns and Categories

We focused on the Modality, Organ System, Plane, and Abnormality categories of questions from Med-VQA-2019 datasets. The category list and related question frequency in the VQA-Med-2019 dataset are given in Table 1.

#### 4.4.1.1. Modality

-Closed, Yes/No, and WH modalities. Examples:

- Did the patient get GI contrast?

- What method was employed to capture this image?

- What does this image's MR weighting look like?

- Is this a flair, t1 weighted, or t2 weighted image?

#### 4.4.1.2. Organ System

WH questions for the organ system. Examples:

- What organ does this MRI primarily show?

- Which organ system does this X-ray show?

### 4.4.1.3. Plane

WH queries. Examples:

- What kind of plane is this mammogram performed in?

- What is this MRI's plane?

### 4.4.1.4. Abnormality

what and Yes/No questions. Examples:

-Does this image appear to be normal?

- What is the image's main abnormality?

- Does this imaging of the gastrointestinal system show any anomalies?

- What about this ultrasound is most concerning?

### 4.4.1.5. The Most Common Answers for Each Category in the VQA-Med-2019

*Table 7:The Most Common Answers for Each Category in the VQA-Med-2019 Training Set.*

| Category | No, of the most frequent questions |
|---|---|
| Modality | An angiography (78), yes (552) no (554), flair (53), mr-t2 weighted (56),ct w/contrast cta-ct angiography,(iv) (50), xr-plain film (456), t2 (217), t1 (137), contrast (107), noncontrast (102), us-ultrasound (183), and mr-flair (84). |
| Organ System | skull and its contents (1216), gastrointestinal (352), musculoskeletal (436), genitourinary (214), sinuses, face, |

|  |  |
| --- | --- |
|  | and neck (191), lymphatic and vascular (122), breast (65), and heart and major vessels (120), |
| Plane | sagittal (478), Axial (1558), coronal (389), pa (92), ap (197), lateral (151), transverse (76), frontal (120), and oblique (50) |
| Abnormality | no (48), yes (62), meningioma (30), pulmonary embolism (16), arteriovenous malformation (avm) (14), schwannoma (13), cerebral abscess, brain (12), fibrous dysplasia (12), ependymoma (12), diverticulitis (11), langerhan cell histiocytosis (11), multiple sclerosis (12), sarcoidosis (11), acute appendicitis (14), tuberous sclerosis (13), glioblastoma multiforme (28), arachnoid cyst (13), |

### 4.4.2. Training and Validation sets for Med-VQA-2019

3,200 images and 12,792 question-answer (QA) pairs contain 3–4 questions for each image making up the training set. The most frequent responses for each category are shown in Table 1. The validation set consists of 2,000 QA pairings and 500 diagnostic images.

### 4.4.3. Test Sets for Med-VQA-2019

The test answers were manually double-validated by a physician and a radiologist. A total of 33 answers were changed, either by (I) introducing an optional component (8 answers), (ii) adding more potential responses (10 answers), or (iii) amending the automated response. 15 answers, or 3% of the total test answers, were modified. Abnormality (8/125), Plane (1/125), and Organ (6/125) are the categories to which the corrected responses belong. For questions involving abnormalities, the adjustment generally involved revising the diagnosis that is implied by the issue visible in the image. The error rate in the training and validation sets, which were produced with the same automatic data production technique, should be comparable. The tested set has 500 questions and 500 medical images.

### 4.5 Med-VQA-2020

The training, validation, and test sets were automatically created by using several filters to choose relevant images and related annotations, and (ii) The questions and their answers are generated by establishing patterns. The test set was manually checked by two medical professionals. The category list and related question frequency in the VQA-Med-2019 dataset are given in Table 2.

### 4.5.1. Question Patterns and Categories

We focused on the Modality, Organ System, Plane, and Abnormality categories of questions from Med-VQA-2020 datasets.

### 4.5.1.1. Abnormality

Medical problems (also with their frequency) examples for the Med-VQA-2020 data:

– acute appendicitis (109)

– pulmonary embolism (114)

– angiomyolipoma (68)

– lung adenocarcinoma (60)

– osteochondroma (63)

– sarcoidosis (58)

*Table 8: The Most Common Answers for Each Category in the VQA-Med-2020 Training Set for abnormality.*

| Category | No, of the most frequent questions |
| --- | --- |
| Abnormality | no (48), yes (62), meningioma (30), pulmonary embolism (16), arteriovenous malformation (avm) (14), schwannoma (13), cerebral abscess, brain (12), fibrous dysplasia (12), ependymoma(12), diverticulitis (11), langerhan cell histiocytosis (11), multiple sclerosis (12), sarcoidosis (11), acute appendicitis (14), tuberous sclerosis (13), glioblastoma multiforme (28), arachnoid cyst (13), – acute appendicitis (109), pulmonary embolism (114), angiomyolipoma (68), lung adenocarcinoma (60), osteochondroma (63) |

### 4.5.2. Training and Validation sets for Med-VQA-2020

4,000 radiology images and 4,000 question-answer (QA) pairs contain 3-4 questions for each image making up the training set. The most frequent responses for the category were shown in Table 2. The validation set consists of 500 QA pairings and 500 diagnostic images.

### 4.5.3. Test Sets for Med-VQA-2020

The test answers were manually double-validated by a physician and a radiologist. A total of 33 answers were changed, either by (I) introducing an optional component (8 answers), (ii) adding more potential responses (10 answers), or (iii) amending the automated response. 15 answers, or 3% of the total test answers, were modified. Abnormality (8/125) is the category to which the corrected responses belong. For questions involving abnormalities, the adjustment generally involved revising the diagnosis that is implied by the issue visible in the image. The error rate in the training and validation sets, which were produced with the same automatic data production technique, should be comparable. The tested set has 500 questions and 500 medical images.

### 4.6  Evaluation Metrics

### 4.6.1.  BLEU

A common metric for assessing NLP systems that generate language, particularly Natural Language Generation (NLG) systems and machine translation (MT), is called BLEU. It is

a metric for evaluating a sentence that was generated against one that was provided. A score of 1.0 indicates perfect matching, whereas a score of 0.0 indicates perfect mismatching. The use of BLEU as an evaluation metric relies on the expectation that it correlates with and anticipates the real function of these systems, measured either externally (for example, by task performance) or by user satisfaction. BLEU itself simply computes word-based similarity with a gold-standard benchmark text.

From this viewpoint, it is comparable to clinical medicine's use of "surrogate endpoints," such as measuring the effect of an AIDS drug on viral load rather than directly determining if it results in prolonged life. Therefore, using BLEU to examine NLP systems makes sense only when their BLEU scores connect with assessments of the direct effectiveness of the value of NLP systems. For the proposed VQA system BLEU score was computed To determine how closely the system-generated responses matched the actual answers on the Med-VQA-2019 and Med-VQA-2020 datasets. The mathematical equation of BLEU for the VQA system is given below.

$$BLEU = BP \cdot \exp(n = 1Nwn \log pn)$$

Here pn was the geometric average of the improved n-gram accuracy, BP was the brevity penalty, and N-grams have positive weight wn= 1/N, and up to length N = 4.

### 4.6.2. Accuracy

For each model, we determined the accuracy, recall, precision, F1-score, specificity, and AUC value to determine performance. These statistical metrics are based on True Positives (TP), False Positives (FP), True Negatives (TN), and False Negatives (FN). The

ratio of correctly detected answers is represented by TN and TP in this example, whereas FN and FP stand for incorrectly classified answers. The following 4 terms are the measures of other performance metrics.

**True Positives (TP):** TP is a situation where "yes" is correctly predicted and the actual result is also "yes."

**True Negatives (TN):** This occurs when "No" is predicted as "No" and the actual result is "No" as well.

**False Positives (FP):** These occur when "yes" is predicted as a result when a no result occurred.

**False negatives (FN):** This occurs when something is predicted to be negative but turns out to be positive.

The accuracy scores indicate how frequently the models delivered the desired answers correctly.

$$Accuracy = (TP+TN)/(TP+FN+FP+FN)$$

### 4.6.3. FI-Scores Metric

These statistics also have an impact on the (F1) score, which is the harmonic mean of recall and precision. By taking into account the harmonic means of recall and precision, the F1-score combines a classifier's recall and accuracy into a singular statistic.

The F1 score is determined by the weighted average as follows:

$$F1\text{-}Score = 2 \text{ (Precision * Recall) (Precision + Recall)}$$

## 4.7 Trans-MedVQA Performance

Figure 7-8 displays the accuracy and loss graph on Med-VQA-2019, it may be inferred from the discrepancy between training and validation accuracy that the model is somewhat overfitting the training data, a flaw common to deep learning models. Our model has achieved around 79. 6% training accuracy and 74.5% validation accuracy. The training loss decreased from 8.6 to 0.44 on the loss graph (Figure. 8), whereas the validation loss decreased from 9.3 to 0.25 after 40 iterations.



*Figure 7: Accuracy graph on Med-VQA-2019*

*Figure 8: Loss graph of Med-VQA-2019*

As can be seen from the accuracy and loss graph on Med-VQA-2020 in Figure. 8-9, the module also experienced overfitting since training accuracy differs from validation accuracy but was less than compared to performance on Med-VQA-2019, which stayed lower than training accuracy. Training and validation accuracy was at 64.1% and 62.1%, respectively, at the conclusion of the previous epoch. Validation loss (Figure. 9) started at 13.5 and decreased to around 0.4 at the last epoch while training loss started at 0.13,4 and decreased to 0.33. At one epoch, validation loss became smaller than training loss, although it remained significant for the majority of the training period.

*Figure 9: Accuracy graph of Med-VQA-2020*



*Figure 10: Loss graph on Med-VQA-2020*

## 4.8 Performance Comparison with state-of-the-art methods

In this subsection, We compare our model's performance with other state-of-the-art VQA methods, such as (MLM) (Khare et al., n.d.), CGMVQA (Ren et al., 2020), BERT (Zhou et al., n.d.), IF-1C (Kornuta et al., n.d.-a), VisualBert ResMLP (Seenivasan et al., 2022), SFN (Kornuta et al., n.d.-a) and CGMVQA Ens. (Ren et al., 2020) on Med-VQA-19, while DenseNet-121 (Liao et al., 2020), VGG-16 with BioBERT (Jung et al., n.d.), NLM (Notes) & 2020, n.d.), LSTM-Encoder-Decoder (Verma et al., 2020), VGG16+GRU+seq2seq (S. Liu et al., n.d.) and Gated Linear Unit (GLU) (Umada et al., n.d.) on Med-VQA-20 dataset.

Our models aimed to improve these state-of-the-art methods for enhancing the accuracy and efficiency of medical VQA. For this comparison, our best performing model, Trilinear Interaction Transformer, is compared with the other state-of-the-art model's performance, as shown below in Table. 5. All the models are pr- trained on the PathVQA and VQARAd datasets, and the results of other models are taken from their respective research articles. Our Trilinear Interaction Transformer outperforms the other research in terms of accuracy, F1-score, blue, and accuracy. What did our model achieved? accuracy? precision? recall? F1-score? blue and? meteor. While the other model's performance accuracy is shown below in Table 9.

*Table 9:Performance comparison with previous state-of-the-art methods.*

| Dataset | Methods | F1-Scores | Blue | Acc. |
|---------|---------|-----------|------|------|
| **VQA-19** | Mask Language Modeling (Khare et al., n.d.) | 0.686 | 0.675 | 0.690 |
| | CGMVQA (Ren et al., 2020) | 61.612 | 0.655 | 0.640 |
| | BERT (Zhou et al., n.d.) | 0.594 | 0.638 | 0.601 |
| | IF-1C (Kornuta et al., n.d.-a) | 0.54 | 56.21 | 0.545 |
| | VisualBert ResMLP (Seenivasan et al., 2022) | 0.601 | 61.21 | 0.736 |
| | SFN (Kornuta et al., n.d.-b) | 0.717 | | 0.710 |
| | CGMVQA Ens. (Ren et al., 2020) | 0.637 | 65.9 | 0.640 |
| | **Our Tri–Transformer-VQA** | **0.724** | **0.784** | **0.746** |

| | | | | |
|---|---|---|---|---|
| | DenseNet-121 (Liao et al., 2020) | 0.442 | 0.54 | 0.49 |
| | VGG-16 with BioBERT (Jung et al., n.d.) | 0.532 | 0.550 | 0.562 |
| **VQA-20** | NLM (Notes) & 2020, n.d.) | 0.425 | 0.448 | 0.405 |
| | LSTM-Encoder-Decoder (Verma et al., 2020) | 0.426 | 0.4358 | 37.8 |
| | VGG16+GRU+seq2seq (S. Liu et al., n.d.) | 0.3914 | 0.410` | 0.375 |
| | Gated Linear Unit (GLU) (Umada et al., n.d.) | 0.6536 | 0.31 | 0.357 |
| | **Our Tri–Transformer-VQA** | **0.624** | **0.639** | **0.621** |

**CHAPTER 5.**

**DISCUSSION**

## 5.1 Discussion of Results

Nearly all human body organs are represented in the VQA data set, along with every type of medical imaging. The dataset sets general questions addressing what doctors should be aware of while interpreting medical images. The majority of deep learning models can simply produce classifications or answers.

We design a Multimodal representation interaction Transformers model that is performed with two transformers having multiple layers where one transformer is for inter-modality input representation and the other Transformer is for intra-modality input representations, to answer questionnaires on medical imagining, their interactions, and surgical procedures based on our two-novel medicinal VQA datasets evolved from two public datasets. To perform classification and both closed and open sentence-based answering, vision-text attention-based transformer models are employed. Trilinear Interaction Attention (TrI-Att) fusion Transformer for the Inter-modality representations interaction, while self-attention for the intra-modality representation interaction (Yu et al., 2019).

A BERT transformer encoder model is used to extract the text embedding from the question and answer, with lesser model parameters that marginally outperform the base

vision-text attention encoder model by incorporating a cross-token sub-module, while EfficientNet is used as the feature extract to get the image embedding. The influence of the number of input image patches and the inclusion of temporal visual features on the model's performance is also reported. Finally, representations obtained from a feed-forward layer are fed to the answer predictor for producing the answer with the answer predictor in an output layer.

While our Medical-VQA system task answers complex questions, from the application standpoint, it unfolds the possibility of incorporating open-ended questions where the model could be trained to answer medical-specific complex questionnaires. Because there were only a limited number of possible responses for each of these question categories, our VQA system did best when answering questions about modality and abnormity, then questions about planes and organs. from a model's perspective. The model has higher accuracy on the elementary questions and as well as can correctly answer complex questions, so it could be used for assisting in teaching beginning medical students or giving the answers to the elementary questions from the patients. Our system can also be beneficial for hospitals because medical imaging is the basis of modern-day medical diagnosis, and hospitals have specialist doctors that see the medical image and determine the disease. Expanding the amount of data can make the model perform better.

## 5.2  Performance discussion of methods

Our model aimed to improve state-of-the-art methods for enhancing the accuracy and efficiency of medical VQA systems. For this discussion, our best performing model, Trilinear Interaction Transformer, is compared with the other model's performance, as shown above in Table. 5.

### 5.2.1.  Med-VQA-19 performance discussion

Our Trilinear Interaction Transformer outperforms the other research in terms of accuracy for the Med-VQA-19 dataset. The results given by Masked Language Modeling (MLM) and CGMVQA were 0.690 and 0.640 with accuracies because these systems only answered the less complex closed-ended questions in terms of modality category. But in the modality category with the largest number of classes, our model has an improvement rate of 6% in accuracy compared to the baseline MLM and CGMVQA methods.

### 5.2.2.  Med-VQA-20 performance discussion

Our Trilinear Interaction Transformer outperforms the other research in terms of accuracy for the Med-VQA-20 dataset. The results given by DenseNet-121 and LSTM-Encoder-Decoder were 0.49 and 0.37 because Open-ended abnormality is difficult to achieve high accuracy. The answers predicted by these methods in the generative mode were phrases

such as ''glioblastoma multiforme''. These words were related to the type of images and the word frequency in the training set. So, our model has an improvement rate of 13% in accuracy compared to the baseline DenseNet-121 and LSTM-Encoder-Decoder methods.

### 5.2.3. Med-VQA-19 and Med-VQA-20 performance discussion

As our Trilinear Interaction Transformer model outperforms the previous methods on both VQA-19 and VQA-20 datasets. For this discussion, on comparison between our model's results, our best method is performing on VQA-19 data with 0.746 accuracy, as compared to VQA-20 with 0.621 accuracy. Because VQA-19 data has less closed-ended question answer categories and with question variations and did not answer the complex questions as VQA0-20. So, our model is 10% less accurate for VQA-20 in terms of accuracy. Because VQA-20 data has a large number of open-ended questions answers it makes them complex to predict an answer for more complex questions correctly with a high accuracy rate.

# CHAPTER 6.

# CONCLUSION, IMPLICATIONS, AND RECOMMENDATIONS

## 6.8 Conclusion

Systems that answer visual questions about medical imagery can be quite beneficial for giving doctors a second view. We proposed the Trilinear Interaction Transformer model in this research, which is a multi-head-attention-based multimodal model for VQA systems on medical images. In order to effectively respond to concerns about medical images and the Trilinear Interaction Transformer model was specifically developed for managing them. On two real-world medical VQA datasets, VQA-19 and VQA-20, a detailed quantitative and qualitative investigation of the performance of the Trilinear Interaction model was conducted for the tasks of answer generation and answer classification in order to produce new state-of-the-art results. Our model improves state-of-the-art methods for enhancing the accuracy and efficiency of the medical VQA system.

## 6.9 Recommendations for Future Research

Future studies could consider further investigating Multimodal fusion techniques as medical diagnosis often involves multiple types of data, such as images and text. Future studies could explore ways to fuse these different modalities within a deep learning model to improve diagnostic accuracy. Also, another possible avenue for research could be Explainable AI, medical professionals need to understand how diagnostic decisions are

being made by computer-aided systems. Hence, developing explainable AI techniques to interpret and visualize the deep learning model's internal workings can improve the trust and adoption of these systems in clinical practice. Also, clinical validation future studies could do to improve our results as it is essential to evaluate the performance of these models using real-world clinical data. Future research can focus on testing the generalizability of the deep learning models on large datasets with diverse patient populations and medical conditions. Future research also needs to be considered the privacy and security aspects, as medical data is sensitive and needs to be protected. the researcher should focus on developing secure and private ways to use deep learning models to aid medical diagnosis. This could involve developing methods to encrypt the medical data before it is processed by the model or using techniques such as federated learning that keep the data on local servers and minimize the amount of data shared between different sites.

## 6.10    Business Benefits of this Research

The research on "Computer-Aided Medical Diagnostic System on Visual Question-Answer Using Deep Learning" may help businesses in several ways.

1. Improved Diagnostic Accuracy: Using deep learning to develop a computer-aided medical diagnostic system can improve the accuracy of medical diagnosis. This may be

extremely advantageous for healthcare organizations since it lowers the likelihood of misdiagnosis and improves patient outcomes.

2. Increased Efficiency and Productivity: Businesses may enhance efficiency and productivity by automating the diagnosis process with deep learning algorithms. The technology can help healthcare workers analyze visual data and answer diagnostic queries more rapidly, allowing them to focus on other vital activities.

3. Cost Reduction: Automated diagnostic technologies have the potential to cut healthcare expenses. The system can improve resource allocation and minimize dependency on expensive diagnostic procedures by reducing the requirement for manual interpretation and analysis of medical images.

4. Knowledge Sharing and Collaboration: A computer-aided diagnostic system can serve as a centralized platform for medical practitioners to share and collaborate on knowledge. The technology can store and analyze massive amounts of medical picture data, allowing healthcare organizations to tap into collective knowledge for precise diagnosis and treatment decisions.

5. Continuous Learning and Improvement: Deep learning algorithms may learn and improve continually over time. The diagnostic accuracy may be improved further by training the algorithm on a large collection of medical photos and related questions and

answers. This iterative learning process can result in continuous improvements in system performance, which will benefit organizations in the long term.

6. Competitive Advantage: Using deep learning to implement a cutting-edge computer-aided diagnostic system can give a competitive advantage to healthcare firms. Businesses may attract more patients, improve their reputation, and differentiate themselves from competition by providing more accurate and efficient diagnostic capabilities.

7. Augmented Decision-Making: The computer-aided diagnostic system can be a valuable tool for healthcare professionals, helping them make better decisions. Doctors can gain additional insights and perspectives by combining their expertise with the system's ability to analyze visual data and provide relevant answers. This can help with complex cases, improve accuracy, and help physicians make well-informed treatment recommendations.

8. Enhanced Patient Care and Experience: Using a computer-aided diagnostic system can lead to better patient care and experience. Patients can receive timely treatment and intervention by enabling faster and more accurate diagnoses. The system can also provide patients with detailed explanations and visualizations, allowing them to better understand their medical conditions. Increased patient satisfaction, trust, and engagement with the healthcare provider can result from this.

The research being conducted on computer-aided medical diagnostic systems that use deep learning for visual question-answering has the potential to benefit the following business areas:

1. Healthcare business: This research has a significant impact on the healthcare business since computer-aided diagnostics may assist medical personnel in properly and efficiently identifying medical disorders. This can result in better patient outcomes and lower healthcare expenses.

2. Medical Technology Companies: Medical technology companies can use the research findings to build improved diagnostic instruments and improve their present products. Deep learning for visual question-answering can aid in the development of more accurate and dependable diagnostic systems.

3. Data Analytics Companies: The research findings may be valuable for data analytics firms that provide services to the healthcare business. These businesses may use deep learning algorithms to analyze medical photos and provide insights into patient diagnosis.

4. Insurance companies: Insurance companies can profit from enhanced medical diagnosis accuracy. More precise diagnoses can assist in avoiding misdiagnoses and needless treatments, resulting in lower costs for insurance providers and patients.

The CEO of a hospital can benefit from the implementation of a computer-aided medical diagnostic system based on visual question-answering using deep learning in several ways:

1. Improved Diagnostic Accuracy: By employing deep learning algorithms to scan medical pictures and answer precise questions about them, the system may considerably enhance diagnostic accuracy. This can result in more accurate and trustworthy diagnoses, lowering the chance of misdiagnosis and improving patient outcomes. The CEO may be certain that the hospital is offering high-quality diagnostic services, which boosts the institution's reputation.

2. Improved Operational Efficiency: Using a computer-aided diagnostic system helps expedite the diagnosis process and enhance hospital operational efficiency. The technology can automate picture processing and deliver quick responses to inquiries, saving time on manual interpretation. This enables healthcare practitioners to make faster judgments, prioritize essential situations, and better allocate resources. The CEO may see improved patient flow, shorter wait times, and higher patient throughput, which leads to greater operational performance and cost savings.

3. Cost Savings: The computer-aided diagnostic system can help the hospital save money. Healthcare resources may be used more efficiently by minimizing unneeded diagnostic

tests and procedures. The system's precise analysis and rapid responses can assist in avoiding unnecessary testing, which can be costly and time-consuming. Furthermore, by lowering misdiagnosis rates, the hospital can avoid costly errors like unneeded treatments or prescription errors. The CEO can anticipate better financial results and resource optimization.

4. Technological Advancement and Innovation: The hospital's commitment to technological advancement and innovation in healthcare is demonstrated by the installation of a cutting-edge computer-aided diagnostic system. This has the potential to attract top talent, such as skilled medical professionals and researchers interested in working with cutting-edge technology. The CEO can establish the hospital as an industry leader, attracting partnerships, collaborations, and funding.

5. Patient Satisfaction and Trust: Computer-aided diagnostic systems can help increase patient satisfaction and trust. Patients can receive appropriate treatment plans faster if accurate diagnoses are provided in a timely manner, leading to improved health outcomes. Patients will appreciate the hospital's dedication to utilizing cutting-edge technologies for their benefit. Patient loyalty, positive word-of-mouth recommendations, and, ultimately, business growth for the hospital can all result from increased patient satisfaction and trust.

6. Data-driven Insights and Research: The computer-aided diagnostic system generates a large amount of data, which can be analyzed to gain valuable insights. This information

can be used by the CEO to identify trends, patterns, and areas for improvement in diagnostic practices. The system can help the hospital's research efforts by allowing for the investigation of new diagnostic techniques, treatment approaches, and medical discoveries. This can help to improve the hospital's reputation as a research institution and foster collaborations between academia and industry.

It should be noted that practical implementation and acceptance of such a system in a corporate setting would include careful considerations such as regulatory compliance, data protection, and interaction with current healthcare infrastructure.

Overall, implementing a computer-aided medical diagnostic system based on visual question-answering using deep learning can provide a hospital CEO with benefits such as improved diagnostic accuracy, increased operational efficiency, cost savings, technological advancement, patient satisfaction and trust, and data-driven insights for research. These advantages contribute to improved healthcare delivery, financial performance, and overall hospital success.

## APPENDIX A

The computer-aided medical diagnostic system

Visual Question-answer

Deep Learning

Multimodal Transformer

Artificial intelligence (AI)

Surgical procedures

Computer-Assisted Surgery

Long short-term memory (LSTM)

Annotated data

Transfer learning

Histopathologic diagnosis

Elementary questions

Health-care images

Voice of the Radiologist

Search engines

VQA dataset

VQA-Med

Open-ended questions

Fine-grained dimension

General language database

Natural Language Processing (NLP)

Computer Vision (CV)

Imbalanced Distribution

QA pair

VQA-RAD

VQA-Med2019

VQA-Med2020

Path-VQA

Closed-ended questions

 MODALITY

PLANE

ORGAN (Organ System)

Gated Recurrent Units (GRU)

ABN (Abnormality)

Local features

PRES (Object/Condition Presence)

Self-attention

POS (Positional Reasoning)

ATTRIB

Recurrent Neural Networks (RNNs)

COUNT (Counting)

Abnormality

Positional reasoning

Classification

Multidisciplinary field

Bounding boxes (detection)

Masks (segmentation)

Medical Question Answering

Cost-effective diagnosis

State-of-the-art

VGGNet

MuVAM model

ResNet

Word-to-text (W2T) attention

Grayscale images

VQA-RADPh

RadVisDial

SLAKE

Visual Objects

Semantic labels

Stacked Attention Network (SAN)

Compact Bilinear Pooling (MCB)

Bidirectional Long Short-Term Memory (BiLSTM)

BERT Model

Bilinear attention networks

MedFuseNet

MICCAI (Medical Image Computing and Computer Assisted Intervention)

MLP-based visualBert encoder

ImageCLEF 2019 VQA-Med data

MS COCO

Single Modality

EfficietNet

Answer Predictor

Embedding of images

ConvNet Architecture

Scalability Technique

Trilinear Transformers

PARALING

Guided attention and self-attention blocks

Softmax

Trilinear interaction. Attention

Intra-modalities

Sequential operations

Absolutely and relatively positional information

Multi-head attention

Feed Forward Layer (FFL)

BLEU Metric

Accuracy

Precision and recall

F1-Scores

Fragmentation penalty

MedPix5 datasets

Training and Validation

Test Data

No, of the most frequent questions

Pre-Training

# REFERENCES

Brown, T. B. et al. (no date) 'Language models are few-shot learners',roceedings.neurips.cc. Available at: https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html (Accessed: 9 February 2022).

Chung, J. et al. (2014) 'Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling'. Available at: http://arxiv.org/abs/1412.3555 (Accessed: 9 February 2022).

Deng, J. et al. (no date) 'ImageNet: A large-scale hierarchical image database', ieeexplore.ieee.org. Available at: https://ieeexplore.ieee.org/abstract/document/5206848/ (Accessed: 9 February 2022).

Devlin, J. et al. (no date) 'Bert: Pre-training of deep bidirectional transformers for language understanding', arxiv.org. Available at: https://arxiv.org/abs/1810.04805 (Accessed: 9 February 2022).

Eslami, S., de Melo, G. and Meinel, C. (2021) 'Does CLIP Benefit Visual Question Answering in the Medical Domain as Much as it Does in the General Domain?' Available at: http://arxiv.org/abs/2112.13906 (Accessed: 9 February 2022).

Fukui, A. et al. (no date) 'Multimodal compact bilinear pooling for visual question answering and visual grounding', arxiv.org. Available at: https://arxiv.org/abs/1606.01847 (Accessed: 9 February 2022).

Goodfellow, I., Bengio, Y. and Courville, A. (2016) Deep learning. Available at: https://books.google.com/books?hl=en&lr=&id=omivDQAAQBAJ&oi=fnd&pg

=PR5&dq=Goodfellow,+I.,+Bengio,+Y.+and+Courville,+A.,+2016.+Deep+lear
ning.+MIT+press.&ots=MNQ2dvlEOT&sig=rlg239q_hC3PYUJZkcqeJwFpqk0
(Accessed: 9 February 2022).

Goyal, Y. et al. (no date) 'Making the v in vqa matter: Elevating the role of image
understanding in visual question answering', openaccess.thecvf.com. Available
at: http://openaccess.thecvf.com/content_cvpr_2017/html/Goyal_Making_the_v
_CVPR_2017_paper.html (Accessed: 9 February 2022).

Greff, K. et al. (no date) 'LSTM: A search space odyssey', ieeexplore.ieee.org.
doi: 10.1109/TNNLS.2016.2582924.Abacha, A. et al. (no date) 'Overview of the
vqa-med task at imageclef 2020: Visual question answering and generation in the
medical domain', openreview.net. Available at: https://openreview.net/pdf?id=Y
PNpmlJcBCd (Accessed: 9 February 2022).

Abacha, A. Ben et al. (2018) 'NLM at ImageCLEF 2018 Visual Question
Answering in the Medical Domain.', researchgate.net. Available at: https://www.
researchgate.net/profile/Asma-Ben-Abacha/publication/328491475_NLM_at_
ImageCLEF_2018_Visual_Question_Answering_in_the_Medical_Domain/links
/5bd0afa745851537f598eb15/NLM-at-ImageCLEF-2018-Visual-Question-Ans
wering-in-the-Medical-Domain.pdf (Accessed: 9 February 2022).

Abacha, A. Ben et al. (2019) 'VQA-Med ':, Proceedings of CLEF (Conference
and Labs of the Evaluation Forum) 2019 Working Notes. Available at: https://aro
des.hes-so.ch/record/4214 (Accessed: 9 February 2022). Alsentzer, E. et al. (no
date) 'Publicly available clinical BERT embeddings', arxiv.org. Available at:
https://arxiv.org/abs/1904.03323 (Accessed: 9 February 2022).

Antol, S. et al. (no date) 'Vqa: Visual question answering', openaccess.thecvf.com. Available at: http://openaccess.thecvf.com/content_iccv _2015/html/Antol_VQA_Visual_Question_ICCV_2015_paper.html (Accessed : 9 February 2022).

Kim, J.-H. et al. (no date) 'Multimodal residual learning for visual qa', proceedings.neurips.cc. Available at: https://proceedings.neurips.cc/paper/2016/ hash/9b04d152845ec0a378394003c96da594-Abstract.html (Accessed: 9 February 2022).

Lau, J. et al. (no date) 'A dataset of clinically generated visual questions and answers about radiology images', nature.com. Available at: https://www.nature.c om/articles/sdata2018251 (Accessed: 9 February 2022). Pan, S., data, Q. Y.-I. T. on knowledge and 2009, undefined (no date) 'A survey on transfer learning', ieeexplore.ieee.org. Available at: https://ieeexplore.ieee.org/abstract/document/5 288526/ (Accessed: 9 February 2022).

Simonyan, K. and Zisserman, A. (2015) 'Very deep convolutional networks for large-scale image recognition, 3rd International Conference on Learning Repres entations, ICLR 2015 - Conference Track Proceedings.

Wang, P. et al. (no date) 'Explicit knowledge-based reasoning for visual question answering', peerj.com. Available at: https://peerj.com/articles/cs353/Reference.z ip (Accessed: 9 February 2022).

Yu, Z. et al. (no date) 'Multi-modal factorized bilinear pooling with co-attention learning for visual question answering', openaccess.thecvf.com. Available at: http://openaccess.thecvf.com/content_iccv_2017/html/Yu_Multi-Modal_Factori zed_Bilinear_ICCV_2017_paper.html (Accessed: 9 February 2022).

Yuan, H. et al. (no date) 'Gold nanostars: surfactant-free synthesis, 3D modeling, and two-photon photoluminescence imaging', iopscience.iop.org. Available at: https://iopscience.iop.org/article/10.1088/0957-4484/23/7/075102/meta (Accessed: 9 February 2022).

Zhan, L. et al. (2020) 'Medical visual question answering via conditional reasoning', dl.acm.org, p. 10. doi: 10.1145/3394171.3413761.

Goyal, Y., Khot, T., Summers-Stay, D., Batra, D. and Parikh, D., 2017. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 6904-6913).

Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C.L. and Parikh, D., 2015. Vqa: Visual question answering. In Proceedings of the IEEE international conference on computer vision (pp. 2425-2433).

Goodfellow I, Bengio Y, Courville A. Deep learning. MIT Press; 2016 Nov 10.
Abacha, A.B., Hasan, S.A., Datla, V.V., Liu, J., Demner-Fushman, D. and Müller, H., 2019, September. VQA-Med: Overview of the Medical Visual Question Answering Task at ImageCLEF 2019. In CLEF (Working Notes).

Abacha, A.B., Datla, V.V., Hasan, S.A., Demner-Fushman, D. and Müller, H., 2020, January. Overview of the vqa-med task at imageclef 2020: Visual question answering and generation in the medical domain. In CLEF (Working Notes).

Litjens, G., Kooi, T., Bejnordi, B.E., Setio, A.A.A., Ciompi, F., Ghafoorian, M., Van Der Laak, J.A., Van Ginneken, B. and Sánchez, C.I., 2017. A survey on deep learning in medical image analysis. Medical image analysis, 42, pp.60-88.

Lau, J.J., Gayen, S., Abacha, A.B. and Demner-Fushman, D., 2018. A dataset of clinically generated visual questions and answers about radiology images. Scientific data, 5(1), pp.1-10.

Greff, K., Srivastava, R.K., Koutník, J., Steunebrink, B.R. and Schmidhuber, J., 2016. LSTM: A search space odyssey. IEEE transactions on neural networks and learning systems, 28(10), pp.2222-2232.

Chung, J., Gulcehre, C., Cho, K. and Bengio, Y., 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. arXiv preprint arXiv:1412.3555.

Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A. and Agarwal, S., 2020. Language models are few-shot learners. arXiv preprint arXiv:2005.14165.

Zhan, L.M., Liu, B., Fan, L., Chen, J., and Wu, X.M., 2020, October. Medical visual question answering via conditional reasoning. In Proceedings of the 28th ACM International Conference on Multimedia (pp. 2345-2354).

Abacha, A.B., Gayen, S., Lau, J.J., Rajaraman, S. and Demner-Fushman, D., 2018, September. NLM at ImageCLEF 2018 Visual Question Answering in the Medical Domain. In CLEF (Working Notes).

Tan M, Le Q. Efficientnet: Rethinking model scaling for convolutional neural networks. International conference on machine learning 2019 May 24 (pp. 6105-6114). PMLR.

Gong, H., Huang, R., Chen, G. and Li, G., 2021. Sysu-hcp at vqa-med 2021: Pan, S.J. and Yang, Q., 2009. A survey on transfer learning. IEEE Transactions on Knowledge and data engineering, 22(10), pp.1345-1359.

Litjens, G., Kooi, T., Bejnordi, B.E., Setio, A.A.A., Ciompi, F., Ghafoorian, M., Van Der Laak, J.A., Van Ginneken, B. and Sánchez, C.I., 2017. A survey on deep learning in medical image analysis. Medical image analysis, 42, pp.60-88.

Pan, S.J. and Yang, Q., 2009. A survey on transfer learning. IEEE Transactions on Knowledge and data engineering, 22(10), pp.1345-1359.

Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A. and Agarwal, S., 2020. Language models are few-shot learners. arXiv preprint arXiv:2005.14165.

Kovaleva, O., Sswering. IEEE Access, 6, pp.57923-57932.

Yu, Z., Yu, J., Fan, J. and Tao, D., 2017. Multi-modal factorized bilinear pooling with co-attention learning for visual question answering. In Proceedings of the IEEE international conference on computer vision (pp. 1821-1830)

Devlin, J., Chang, M.W., Lee, K. and Toutanova, K., 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.

Yang, Z., He, X., Gao, J., Deng, L. and Smola, A., 2016. Stacked attention networks for image question answering. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 21-29).

Anderson, P. et al. (no date) 'Bottom-up and top-down attention for image captioning and visual question answering', openaccess.thecvf.com. Available at: http://openaccess.thecvf.com/content_cvpr_2018/html/Anderson_Bottom-Up_ and_Top-Down_CVPR_2018_paper.html (Accessed: 3 March 2022).

Attention is All you Need (no date). Available at: https://proceedings.neurips.cc/p aper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html (Accessed: 7 March 2022).

Devlin, J. et al. (2018) 'BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding', NAACL HLT 2019 - 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings of the Conference, 1, pp. 4171–4186. Available at: https://arxiv.org/abs/1810.04805v2 (Accessed: 3 March 2022).

EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (no date). Available at: http://proceedings.mlr.press/v97/tan19a.html (Accessed: 7 March 2022).

Lin, Z. et al. (2017) 'A Structured Self-attentive Sentence Embedding', 5th International Conference on Learning Representations, ICLR 2017 - Conference Track Proceedings. Available at: https://arxiv.org/abs/1703.03130v1 (Accessed: 7 March 2022).

Vaswani, A. et al. (no date) 'Attention is all you need', proceedings.neurips.cc. Available at: https://proceedings.neurips.cc/paper/7181-attention-is-all-you-need (Accessed: 7 March 2022).

Wu, Q. et al. (2017) 'Visual question answering: A survey of methods and datasets', Computer Vision and Image Understanding, 163, pp. 21–40. doi: 10.1016/J.CVIU.2017.05.001.