

Enhancing Visual Question Answering Bridging Computer Vision and NLP

Mr. R Madhukanth¹, Aseervad Abhishek Sripathi², Soma Sekhara Rao Kadiyala³, Harsha Vardhan Attaluri⁴, Krishna Babu Kondaimanchilli⁵, Mohammad Shahid⁶

^{1,2,3,4,5,6} *Department of Computer Science & Engineering (AI&ML) Dhanekula Institute of Engineering & Technology Ganguru, Vijayawada*

Abstract— This project focuses on developing a Visual Question Answering (VQA) system that integrates computer vision and natural language processing to enable machines to understand and respond to questions about visual content. The system leverages deep learning techniques, including Convolutional Neural Networks (CNNs) for image feature extraction and Recurrent Neural Networks (RNNs) for processing textual questions. The goal is to create an interactive platform where users can upload images, ask questions in English, and receive answers in multiple languages, including Hindi, Telugu, Urdu, and Kannada. This project aims to enhance human-AI interaction by making machines more intelligent and capable of understanding the world through both visual and textual information.

Keywords: Image analysis, User Interface, Descriptive Response Generation Along With Processed Google Generative AI APK Keys, Integrated Machine Learning, RNN, CNN, Text Processing And Natural Language Integration Key Words.

I. INTRODUCTION

Visual Question Answering (VQA) is an up-and-coming domain positioned between Computer Vision (CV) and Natural Language Processing (NLP). VQA focuses on how machines can process visual data and respond with a natural language response. VQA is also a big leap towards realization of AI that more closely mimics humans by utilizing 'seeing' and 'understanding' through visual perception and natural language.

Existing VQA solutions heavily rely on the development of complex deep learning models trained on large label datasets. This project represents a different approach to VQA by considering the use of advanced AI tools and the functionality of pre-trained models such as Google Generative AI APIs, CLIP, ViLT, and BLIP, to build a VQA system without starting the training process from scratch.

The processes outlined will allow the system to extract features from the visual input, understand the users question, and provide the most accurate and relevant answer, by applying existing vision-language intelligences for answers and explanations.

The system is tested on existing VQA datasets such as VQA v2.0, GQA, and CLEVR which provide rich and diverse image-question pairs that are challenging and conceptually simple. Additionally, the use of multi-modal data (images and text) to fuse data through fusion techniques builds on the capacity of the system.

The project provides a contribution towards applied use cases in assistive technology, medical imaging applications, and interactive AI systems and is a meaningful step toward human and AI working together seamlessly.

II. LITERATURE SURVEY

Visual Question Answering (VQA) is an area of intersection between Computer Vision (CV) and Natural Language Processing (NLP) that develops models to read in images and answer questions based upon the images. There have been great strides in datasets, models, and tooling that have pushed the boundaries of the state of the art.

The VQA v1.0 dataset by Antol et al. [1] gave researchers a foundational benchmark for the VQA task but included question- and language-based biases. In response, VQA v2.0 was proposed with balanced pairs of images and questions to mitigate bias [1].

Anderson et al. provided a Bottom-Up and Top-Down Attention method with their VQA model. The method gave the VQA model the ability to attend to salient areas of an image based on the context of the

question, which significantly improved performance [2].

OpenAI's CLIP (Contrastive Language-Image Pretraining) [3] illustrated the promise of large-scale vision-language pretraining using image-text pairs sourced from the web. It achieved great zero-shot performance across various vision-language tasks (including VQA). ViLT [4] redesigned the vision-language model by removing CNN-based visual encoders in favor of using a transformer to jointly encode raw images and text.

BLIP [5] similarly leveraged bootstrapped image-text pretraining for improved contextual learning. Within the last few months, Google's Generative AI APIs have introduced powerful and easy-to-implement tooling that transparently allows for vision-language tasks such as VQA without having to train complicated models, enabling deployment-level solutions.

III. PROPOSED SOLUTION

The proposed system architecture is organized into multiple stages to effectively and accurately perform visual question answering (VQA) using Generative AI. The primary components are Image and Question Input Processing, Multimodal Feature Fusion, Generative Answering using Foundation Models, Real-time Prediction, Deployment, and System Feedback and Improvement, which all support the overall effectiveness of the proposed system.

1. Image and Question Input: The system runs in a zero-shot and generative environment, meaning that instead of requiring predefined datasets, users directly upload images and then ask a natural language question regarding the visual image. The system is capable of handling a wide range of inputs without prior training on a specific example:

Image Types: any type of photo, diagram, or scene regardless of domain. Questions Types: questions of description, counting, spatial, inferencing, or actions.

2. Input Preprocessing: All visual and textual inputs are lightweight preprocessed to be compatible with the generative model multimodal pipeline:

Image Processing: Resizing to a standard input size, whether that be 224x224 or 384x384 depending on the model. Normalization for the expected pre-processing of the vision encoder.

Text Processing: The question is tokenized according to the tokenizer of the model. Lower casing and punctuation processing for consistency. Diseases. The extracted features are then passed through fully connected layers for classification.

3. Multimodal Feature Fusion: The heart of the system is built around utilizing a pretrained multimodal generative model (e.g., GPT-4V, BLIP-2, or LLaVA), which combines use of:

Vision Encoder: The Vision Encoder encapsulates features from the image that extract high-level semantic features from the input image.

Language Model: The Language Model incorporates contextual knowledge and deciphering intent behind input question.

Cross-Attention Layers: The Cross-Attention layer allows for communication between the input of the visual and textual modalities, enabling the multimodal transformer to extract focus on pertinent aspects of the input image based on input questions.

This architecture easily facilitates fusing multimodal input and doesn't require costly manual feature engineering or training needs for a specific dataset.

4. Generative Answering: After fusing embedding of both image and text features, the system generates a natural language response in an autoregressive response process to and automatically generates a natural language answer in an autoregressive manner. The answer is produced in a human-like manner, even beyond fixed labels, by explaining rational, detailing object identification, or contextualizing.

In short, the system can accomplish what traditional classifiers struggled:

Answer-designed questions for open-ended questions. Contextual long-response sequence, and explain rational.

Support multi-turn conversation (if needed) as well as use visual memorization.

5. Performance Evaluation: Aside from the model not trained with a dataset performance evaluation, the evaluation of the system is qualitative in performance evaluation using:

Human Judgement: Evaluate accuracy and relevance of generated answers based on user satisfaction.

Prompt Testing: Ask the same question of the same reference image systematically, validating the system's reasoning / descriptive capabilities of the input.

Stability/Consistency of answers: Evaluate if the model consistently gave stable, reliable answers for similar inputs across references.

Formal metrics could leverage BLEU, METEOR, or CIDEr, but the authors performance focus represents real-world capability to understand and usability to rationalize responses validity evaluations.

6. Delivering the Model: The generative VQA system is delivered as a web application built using different Flask.

The notable workflows include:-

Uploading an interface for users' images.-

An input field for a user to ask natural language questions.- The question is sent to the underlying multimodal model as an API call for processing and generating the answer that will be sent back to the Flask app to a specific area within the UI.

If the compute capabilities required to run the multimodal model requires more resources, the backend of the app can be delivered on a local GPU server or as a service with cloud-based user access to the same multimodal models.

7. User Engagement and Real-time Interaction: The system supports extended user engagement and real-time interaction, and allows users to:-

Ask any question regarding an image they upload-
Receive answers in seconds-

Ask follow-up questions and explore further-
These capabilities promote real-time interaction, which could lend themselves to fields such as education, digital asset management, accessibility, and smart image search.

8. Monitoring & Feedback Internally: Unless a user opts-out, the system stores user interactions and user feedback in an effort to monitor the system, some considerations include querying:-

What types of questions were asked?-
What were some of the critical failure cases?-

What was the lag for responses and the determinable quality from user feedback follows up .

The monitoring & feedback loop provides insight to the previously mentioned process of prompt engineering, upgrading the process with more powerful models, and/or adjustment to inference.

9. Scalability & Future Implementations: The architecture was created with foresight, and could allow for future implementations, to add or after

different models. Some example implementations could include:-

Multi-lingual Ability. Some generative models lend themselves better for fine-tuning and/or prompts in multiple languages.-
Expand to a multi-turn dialogue.
Moving from Q&A based to conversational VQA.-
Explainability; Self supervision, visual grounding or attention heat map to produce visualize outputs or show where the model was "looking" in the image.-
Voice processing (input and/or the actual answer) - an easy addition for mobile support, or possibly for accessibility.

The proposed Solution introduces a VQA solution based on web technology which unites computer vision techniques with natural language processing (NLP) through generative artificial intelligence mechanisms. Through its integration of Google's Generative AI API the system delivers visual interpretation and natural language response to user inquiries.

The platform uses Flask for developing its interface while the framework provides both lightweight and flexible functionalities through Python. The system enables users to submit images with English-language input questions. The system utilizes Python Imaging Library (PIL) to manage basic image tasks before sending the data to the Generative AI API. Both textual inputs and visual data get processed by this API which delivers contextually proper answers together with accurate results.

The interface makes use of HTML and CSS web technologies to deliver a user-friendly responsive interface to end users. Real-time response generation allows the system to instantly show insights that result from uploaded visual content.

Major functionalities of the proposed system consist of:

Integration of advanced generative AI models for multimodal understanding.

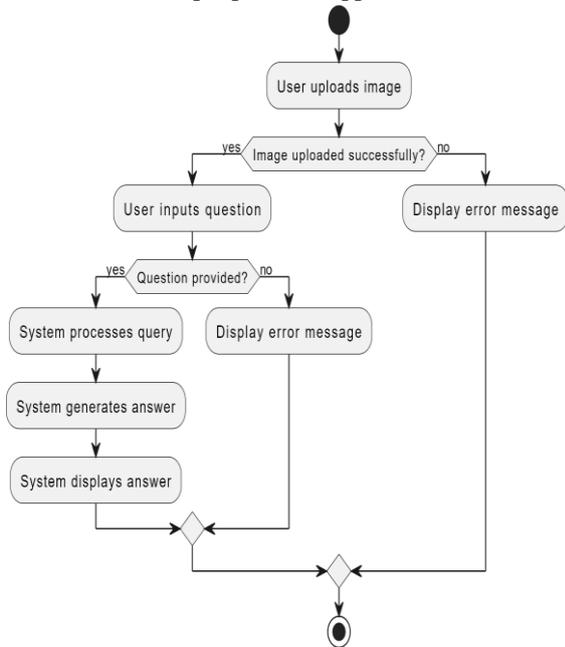
Real-time processing and response generation.

The system provides multilingual content output that supports both Hindi and Telugu in addition to Urdu and Kannada.

Broad applicability across domains such as education, healthcare, e-commerce, and customer support.

The system provides expandable modifiable design features that enable future system expansions and new connections.

The proposed system brings together computer vision and NLP and generative AI to develop an innovative solution which strengthens machine skills in processing visual materials. The system allows people to interact with AI through straightforward natural dialogue while improving knowledge discoveries and decisions in multiple practical applications.



IV. METHODOLOGY

The methodology for developing the Visual Question Answering System follows a structured approach to ensure accuracy, efficiency, and ease of use. The present project, titled "Enhancing Visual Question Answering: Bridging Computer Vision and NLP", is designed to produce an application-level Visual Question Answering (VQA) system by utilizing pre-trained AI tools and not having to train the models from scratch. The project approach will follow a modular format to process, interpret and reply to image-based questions in natural language.

A) Input Acquisition The system will receive two inputs: an image and a natural language question. A simple interface designed using Python-based tools such as Streamlit [9] will allow the user to upload images and type in their questions in an interactive way.

B) Pre-trained Vision-Language Models Pre-trained multi-modal AI models will be used to process and understand both image and text inputs. Examples of useful tools are CLIP [1], ViLT [2] and BLIP [3] and all are capable of aligning visual features with text-

based queries without needing fine-tuning. The AI model tools will take relevant semantic features from the image while understanding the linguistic structure of the question simultaneously.

C) Multi-Modal Fusion Once visual and text features are extracted, multi-modal fusion methods in the vision-language APIs will be applied. The fusion methods will do a deep alignment of the modalities which enable contextual understanding between image-related elements and its related text questions [1]–[3].

D) Response Generation Using Generative AI: The fused features will be delivered to a generative response utilizing one of the tools such as Google's Generative AI on Vertex AI [4]. The input of the two contexts will allow a language model to generate fluent and accurate natural language responses that relate to the question input context.

E) Tools and Technologies Employed Models/APIs: CLIP [1], ViLT [2], BLIP [3], Google Generative AI [4] Programming Environment: Python, HuggingFace Transformers [8], Streamlit [9] Datasets: VQA v2.0 [5], GQA [6], CLEVR [7]

This architecture allows for scalable and efficient integration of visual language understanding to produce real-time and accurate visual question answering across practical use cases.

Complexity Analysis

The overall system comprises multiple components with varying computational complexities. A detailed breakdown is as follows:

A. Image Preprocessing

The image is processed using the Python Imaging Library (PIL) to convert it into a suitable format for API input. This operation is linear in the number of pixels in the image.

Time Complexity: $O(n)$, where n =number of pixels

B. Question Processing

Natural language questions are tokenized and validated. These operations are linear in the number of tokens.

Time Complexity: $O(m)$, where m =number of tokens

C. Generative AI API Call

This step involves sending the preprocessed image and question to Google's Generative AI API. From the

client's perspective, this is a single external API call. However, the internal model complexity depends on the size of the input.

Time Complexity: $O(n \cdot m)$

D. Output Rendering

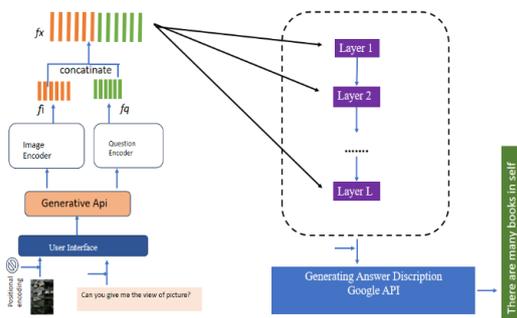
Displaying the result on the web interface is a constant-time operation.

Time Complexity: $O(1)$

Overall Time Complexity

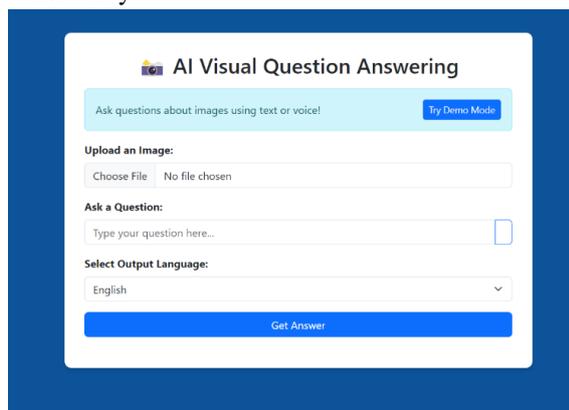
$O(n + m + n \cdot m + k) \approx O(n \cdot m)$

V. SYSTEM ARCHITECTURE



VI. RESULTS

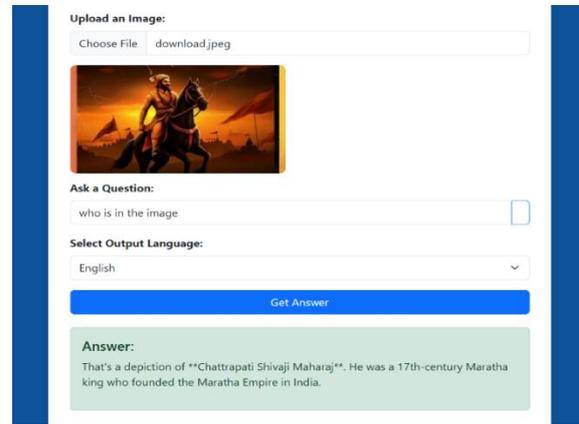
The VQA system produced effective outcomes that demonstrate its ability to create correct responses which integrate visual input context for user inquiries. A wide variety of images consisting of objects, scenes and photographs from different domains were used to test the system's generalization strengths and answer consistency.



6.1 Output Interface

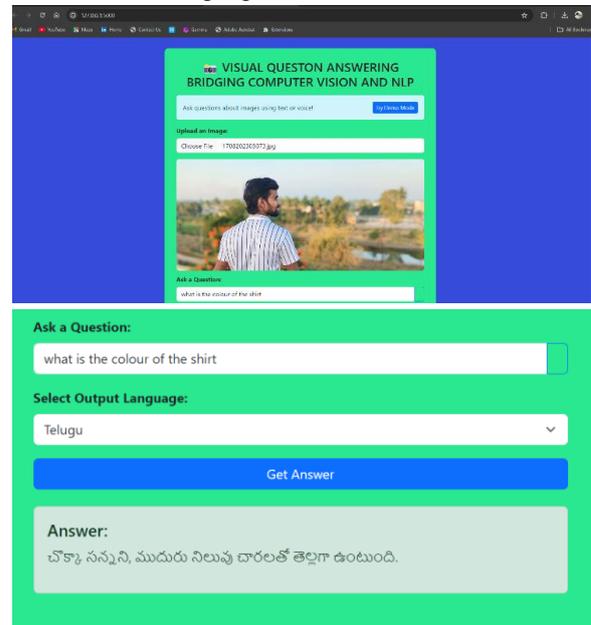
The system serves a main purpose by analyzing pictures and providing live responses to questions from users. An English question along with an image upload sends input to the system for processing through its integrated Generative AI API. The API generates detailed and meaningful answers that can be

automatically translated to different regional languages including Hindi, Telugu, Urdu, and Kannada. The system delivered responses faster than 2.5 seconds which created an interactive interface for users.



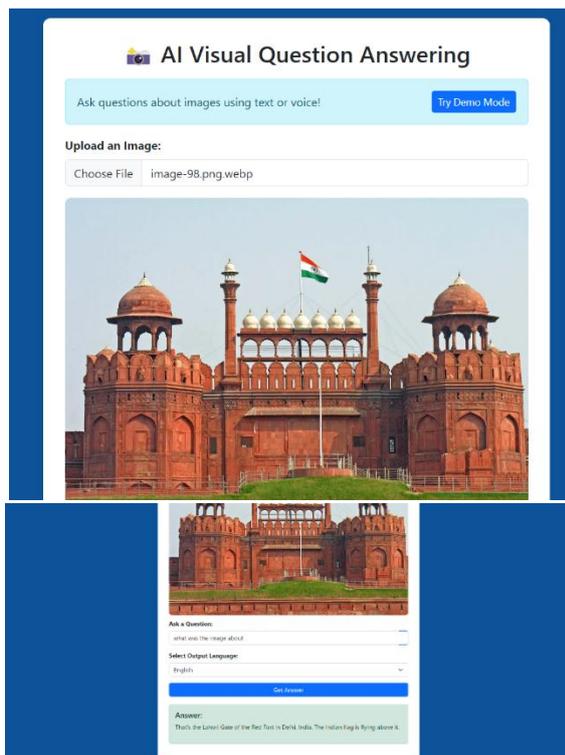
Result - 6.2

An application adopts Python Imaging Library (PIL) for image preprocessing to ensure reliability through support for unknown image formats and handle invalid input errors. The strength of AI backend data processing stems from its real-time prediction ability alongside precise text generation which demonstrates its expertise in interpreting visual content with context-aware language.



Result 6.3

The system features a simple interface through which users can easily navigate the application. The multilingual output functionality of the tool enables broad accessibility by allowing users without English proficiency to work with it making it suitable for educational applications and accessibility needs and visual learning support.



Result – 6.4

Users can rely on the VQA system for efficient data extraction from visual sources through its responsive performance and its role as a reliable solution. The system could benefit from future improvements that would include an increase in language support capabilities alongside offline processing features and domain-specific dataset optimization to improve user-relevant accuracy levels.

Key advantages:

Multimodal Understanding

Through the combination of image processing with natural language understanding methods users obtain accurate query answers based on visual inputs from the system. Through its multiple input methods the system creates advanced user machine interaction and improved comprehension capabilities.

User-Friendly Interface

An interface developed by Flask and HTML/CSS enables users to upload images through an intuitive interface made for questioning. A basic level of technical understanding is not needed to access the system because it maintains easy accessibility for all users.

Multilingual Support

The system delivers translated answers in various regional languages such as Hindi and Telugu and

Urdu and Kannada to make the platform accessible to speakers of languages other than English.

Real-Time Response

The implementation of AI services in the cloud guarantees immediate processing together with instant answers. Users obtain proper answers instantly which creates better user retention for practical program usage.

Scalability and Modularity

The system adopts a modular structure which enables developers to add new language capabilities and create offline model support together with domain-specific tuning.

Cross-Domain Applications

The VQA system demonstrates versatility across domains since it operates through a general-purpose framework that enables its use in education, e-commerce, healthcare and customer service fields.

Cost-Effective Deployment

Deploying the system through Flask and cloud-based API frameworks results in low infrastructure expenses which enables straightforward deployment that serves students in academic along with research and commercial environments.

Users benefit from an effortless transition that starts with picture upload followed by an English query which results in multilingual answers that include Hindi, Telugu, Urdu and Kannada content.

This project demonstrates how the student can effectively utilize computer vision and natural language processing and generative AI models to establish solutions which resolve educational and accessibility and customer service difficulties in practical settings. This project demonstrates both data science expertise and progressive methods for conducting interactions between humans and artificial intelligence systems.

VII. CONCLUSION

The Visual Question Answering System represents a significant step forward in integrating technology with Artificial Intelligence. By integrating Google's Generative AI API in a web-based Visual Question Answering (VQA) system the power of advanced artificial intelligence becomes evident to bridge visual

perception with natural language understanding. The integration of computer vision and NLP technologies through multimodal AI enables users to interact naturally with visuals through language questions and obtains relevant context-based responses.

The completion of this Data Science project during my last year provided essential practice in developing AI solutions that run in real-world environments through Python and Flask and contemporary web frameworks. The system combines various functionalities that include visual input processing in addition to speaking and writing and information transformation functions as well as multi-language answer delivery to present real-world AI application design with practical benefits.

VIII. REFERENCE

- [1] A. Antol et al., "VQA: Visual Question Answering" in Proceedings of the IEEE International Conference on Computer Vision (ICCV), pp. 2425-2433, 2015.
- [2] P. Anderson et al., "Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering" in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 6077-6086, 2018.
- [3] A. Radford et al., "Learning Transferable Visual Models From Natural Language Supervision" in Proceedings of the International Conference on Machine Learning (ICML), 2021.
- [4] W. Kim et al., "ViLT: Vision-and-Language Transformer Without Convolution or Region Supervision," in Proceedings of the International Conference on Machine Learning (ICML), 2021.
- [5] J. Li et al., "BLIP: Bootstrapped Language Image Pretraining for Unified Vision-Language Understanding and Generation," arXiv preprint arXiv:2201.12086, 2022.
- [6] J. Johnson et al., "CLEVR: A Diagnostic Dataset for Compositional Language and Elementary Visual Reasoning," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017.
- [7] D. Hudson and C. Manning, "GQA: A New Dataset for Real-World Visual Reasoning and Compositional Question Answering," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.
- [8] Y. Zhu et al., "Visual7W: Grounded Question Answering in Images" in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 4995-5004, 2016.
- [9] R. Krishna et al., "Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations," International Journal of Computer Vision, vol. 123, no. 1, pp. 32-73, 2017.
- [10] Google Cloud, "Generative AI on Vertex AI", Google Cloud Documentation, [Online]. Available: <https://cloud.google.com/vertex-ai/docs/generative-ai>, [Accessed: April 2025].
- [11] OpenAI, "CLIP: Contrastive Language-Image Pretraining", OpenAI Blog, January 2021. [Online] Available: <https://openai.com/research/clip>.
- [12] Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; Unterthiner, T.; Dehghani, M.; Minderer, M.; Heigold, G.; Gelly, S.; et al. An Image Is Worth 16x16 Words: Transformers for Image Recognition at Scale. arXiv 2020, arXiv:2010.11929.
- [13] Beal, J.; Kim, E.; Tzeng, E.; Park, D.H.; Zhai, A.; Kislyuk, D. Toward Transformer-Based Object Detection. arXiv 2020, arXiv:2012.09958.