International Journal of Scientific Research in Engineering and Management (IJSREM)

Volume: 09 Issue: 05 | May - 2025

SJIF Rating: 8.586

Conversational Image Recognition Chatbot PROF. JACOB AUGUSTINE¹, T ABHINAV KISHAN², K SAI KARTHIK³, T YESHWANTHTEJA⁴, BANDI ANIL KUMAR⁵

¹Professor in Computer science and Engineering & presidency university, Bengaluru 2345 Student in Information Science and Technology & presidency university, Bengaluru ***<u>*</u>***

Abstract - This research introduces a Conversational Image Recognition Chatbot powered by generative AI, designed to enable users to interact naturally with both text and visual inputs. The system combines intelligent dialogue with image understanding, allowing users to upload images, ask context-specific questions, and receive meaningful responses in real time. At its core is Google's Gemini

2.0 Flash model, which supports multimodal processing and delivers fast, focused replies based on both textual prompts and visual content. The application is built using the Flask framework on the backend, with a clean, responsive frontend powered by HTML templates and Bootstrap. It handles image uploads using the Python Imaging Library (PIL) and supports RESTful API endpoints for seamless communication between the client and server. Adjustable AI parameters offer flexibility in response generation, while error handling ensures stable, user-friendly performance. This chatbot bridges the gap between image recognition and natural language conversation, offering a unified platform for educational, creative, and productivity- based use cases. The project demonstrates how generative AI can be effectively applied to enhance human-computer interaction by making both visual and verbal inputs part of a single conversational experience.

Keywords: Conversational AI, Image recognition, Generative AI, Multimodal chatbot, Flask backend, Gemini model, Responsive web interface

INTRODUCTION

The AI-Powered Image and Text Analysis System represents a leap forward in how we interact with technology, blending advanced artificial intelligence with real-world applications. Built using the Flask web framework and integrated with Google's Gemini AI model, this system is designed to handle both text conversations and image analysis, offering a flexible and user-friendly platform.

At the heart of this system is the Gemini-2.0-flash model, carefully fine-tuned with settings (temperature: 0.7,

top p: 0.8) to produce balanced and accurate results. The platform is built with a dual-interface design: one side is a conversational chatbot for text-based communication, and the other is a powerful image analysis tool that can process and describe images.

The system follows a clean, RESTful architecture, with separate endpoints for handling text and image requests. It also includes robust error handling and response management to ensure smooth performance. For images, the system supports a variety of formats using the Python Imaging Library (PIL), while the text component is designed to generate thoughtful, context- aware responses.

In essence, this system takes cutting-edge AI and makes it accessible, allowing users to engage with powerful tools for text and image analysis in a straightforward, easy-to-use interface. It's a perfect example of how modern AI can simplify complex tasks and enhance everyday interactions.

1. LITERATURE SURVEY

Recent advancements in AI-powered image and text analysis systems have led to significant improvements in various areas, particularly in natural language processing, computer vision, web-based AI applications, and user interface design. In the field of natural language processing, large language models, like Google's Gemini AI, have greatly enhanced contextual understanding, enabling more coherent and relevant responses. These models also excel at managing multiturn conversations and maintaining context throughout interactions, making them ideal for conversational AI systems. On the other hand, modern computer vision systems have made tremendous strides in object detection, recognition, and scene understanding. They can now describe complex scenes in detail and integrate both text and image data simultaneously, which enhances their functionality. Real-time image processing capabilities have further expanded the potential applications of these systems.



Flask-based AI applications have gained traction due to their lightweight and efficient architecture, allowing for easy integration with AI models. Flask's support for RESTful APIs also makes it an attractive choice for developers, as it enables scalable and flexible deployment of AI-powered systems. Additionally, recent studies have shown that API-first architectures, microservices-based deployments, and cloud-based AI hosting offer efficient solutions for integrating AI services, ensuring smoother operations and improved scalability. User interface design has also evolved, with effective UI patterns focusing on intuitive chat interfaces, real-time response visualization, and accessibility features that enhance user experience. Security and performance have not been overlooked, with emphasis placed on secure API key management, rate limiting, response optimization, and robust error handling to maintain system reliability and stability.

Looking ahead, future trends suggest further improvements in multi-modal capabilities, enabling AI systems to process and integrate various input types more seamlessly. There is also a push for more accurate and context-aware responses, alongside advanced error recovery mechanisms to minimize disruptions in AIdriven applications. These ongoing advancements demonstrate the rapid evolution of AI systems and their growing ability to deliver sophisticated, user-friendly solutions across various industries.

3. PROPOSED METHOD

3.1 System Architecture

The AI-powered Image and Text Analysis System is designed with a modular architecture to efficiently manage its components, ensuring smooth functionality and integration.

3.1.1 Backend Framework

- Flask Web Server: At the heart of this system is the Flask web server, which provides a simple yet powerful backend framework. It ensures fast deployment and easy integration for both text and image-based operations.
- **RESTful API Design**: The system communicates with external components through a RESTful

API, making it easy to process requests for both chat interactions and image analysis tasks.

• Error Handling System: To ensure seamless user experience, the system includes robust error handling to catch any issues during communication or data processing.

3.1.2 AI Integration Layer

- Google Gemini AI Model: The AI model driving the text and image analysis is Google's Gemini- 2.0-flash, known for its powerful language understanding and image processing capabilities.
- Model Configuration:
 - **Temperature**: Set to 0.7 to provide balanced creativity, ensuring outputs that are both relevant and creative.
 - **Top-p**: A setting of 0.8 ensures the model focuses on generating coherent and meaningful responses.
 - **Top-k**: With a value of 40, this parameter controls the selection of tokens, ensuring more accurate and precise responses.
 - **Max Output Tokens**: Limited to 150 tokens to keep responses concise and relevant to the user's needs.

3.1.3 Image Processing Pipeline

- **PIL Integration**: The Python Imaging Library (PIL) is integrated to handle the preprocessing of images. It ensures that images are processed effectively before being passed to the AI model for analysis.
- **Image Data Management**: The system manages image data efficiently, allowing for smooth image conversion and analysis.
- **Multi-format Support**: It supports various image formats, ensuring that the system can handle a wide range of image types for analysis.

3.3 System Flow

1. Request Reception:

• The system begins by receiving incoming requests, validating them, and preprocessing the data before routing it to the appropriate component, either



the text-based chat or the image analysis part of the system.

2. AI Processing:

 Based on the nature of the input (text or image), the appropriate AI model is selected. The system applies the right configurations to the model to generate accurate and relevant responses or image analysis results.

3. Response Handling:

 Once the AI has processed the request, the response is formatted appropriately, error-checked, and sent back to the client. The system ensures that all responses are clear, concise, and reliable.

3.4 Technical Specifications

- Language: Python 3.x
- Web Framework: Flask
- **AI Model**: Gemini-2.0-flash
- Image Processing: PIL/Pillow
- **Response Format**: JSON

3.3 PROJECT WORKFLOW



4. METHODOLOGIES

1. System Overview

This project aims to bring together Google's Gemini AI technology with a user-friendly web interface to enable text and image analysis. Below is a step-by-step breakdown of how the system works and its key components:

2. Core

Components

Frontend Layer:

- A welcoming landing page
- A dynamic, interactive chat interface
- An easy-to-use image upload feature
- A clear, concise response display system

Backend Architecture:

- Powered by Flask, the backend handles all the server-side logic
- Manages how requests are routed and processed
- Ensures smooth and consistent responses

Т



A robust system for handling errors

AI Integration Layer:

- Integrates Google's Gemini AI model to process text and images
- A pipeline to process text and analyze context
- A pipeline to handle image analysis
- Generates responses based on the processed inputs

3. Processing

Pipelines Text

Processing:

- Ensures that the input text is valid and wellformed
- Manages the context for better, more relevant responses
- Formats the output text for easy readability
- Handles any errors in text processing

Image Processing:

- Validates image formats before processing •
- Uses PIL for any necessary image preprocessing
- Analyzes the image with the AI model
- Generates responses based on the analysis

4. System Features

Model

Configuration:

- **Temperature**: Set to 0.7 for a good balance of creativity and precision
- **Top-p**: Set to 0.8 for focused and relevant outputs
- **Top-k**: Set to 40 to control the range of token selections
- Max Output Tokens: Limited to 150 tokens for concise, informative responses

Error Management:

- Validates inputs to avoid processing errors
- Handles any exceptions that might occur during the process

- Verifies that the responses are correct
- Includes fallback mechanisms in case of system failures

Response Handling:

- Uses a clean and structured JSON format for easy parsing
- Provides error status codes when something goes wrong
- Sends success messages when operations are completed successfully
- Optimizes responses for efficiency

5. Security Implementation

- API key management ensures secure access
- Validates each request to prevent malicious activity
- Sanitizes responses to avoid security vulnerabilities
- Logs errors for further review and troubleshooting

6. Performance Considerations

- Caches responses to reduce load times
- Optimizes image sizes for faster processing •
- Implements rate-limiting to control the number of requests
- Manages server resources to ensure smooth • operation under load

5. RESULT

The implementation of the AI-powered image and text analysis system has been highly successful, integrating Google's Gemini AI with a Flask-based web application. The system efficiently processes text and image inputs, generating optimized and context-aware responses with precise image analysis. The AI model has been configured with carefully tuned parameters, achieving balanced outputs, focused responses, and concise text generation. The image processing pipeline, supported by PIL, ensures fast and accurate handling of images. The user interface is intuitive, featuring a responsive landing page, interactive chat functionality, and seamless image upload capabilities. The system's performance is reliable, with quick response times, stable processing even under load, and robust error



handling. Security measures, including secure API key management, request validation, and response sanitization, have been successfully implemented. Overall, the system offers a reliable, user-friendly solution, combining cutting-edge AI capabilities with smooth web-based functionality for real-time text and image analysis.

LANDING PAGE



CHATBOT INTERFACE



IMAGE FOR TESTING



WORKING OF ANALYZED IMAGE



6. CONCLUSION

This project effectively combines advanced AI with web technologies, demonstrating how Google's Gemini AI model can be integrated with Flask to power real-time text and image analysis. With smooth image processing using PIL, strong error handling, and secure API communication, the system provides fast text processing, accurate image descriptions, and reliable performance. The user-friendly interface ensures an intuitive experience, offering clear responses and helpful error messages. Looking ahead, there's room for growth, such as adding more input formats, enhancing customization, and extending model capabilities. In the end, the project delivers on its goals, laying the groundwork for future AI-powered web applications.

REFERENCES

Liu, Z., Lu, Y., & Song, L. (2020) 'Deep learning for conversational image recognition: A review', *IEEE Access*, 8, pp. 48590-48603.

Joulin, A., Grave, E., Mikolov, T., & Ranzato, M. (2017) 'Bag of Tricks for Efficient Text Classification', *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics*, pp. 427–431.

Vinyals, O., & Le, Q. V. (2015) 'A Neural Conversational Model', *Proceedings of the 32nd International Conference on Machine Learning (ICML* 2015), pp. 1–19.

Yin, X., & Liu, Z. (2021) 'Visual Dialog: A Survey of Conversational Image Recognition Systems', *Computer Vision and Image Understanding*, 202, pp. 103140.

T



Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., & Parikh, D. (2015) 'VQA: Visual Question Answering', *Proceedings of the IEEE International Conference on Computer Vision (ICCV 2015)*, pp. 2425–2433.

Das, A., Kottur, S., Gupta, S., Saenko, K., & Batra, D. (2017) 'Visual Dialog', *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2017)*, pp. 326–335.

Zhou, B., & Chen, X. (2020) 'Vision-and-Language Navigation: A Survey', *IEEE Transactions on Neural Networks and Learning Systems*, 31(9), pp. 3275–3294.

Li, X., & Li, J. (2020) 'End-to-End Image Captioning with Transformer Models', *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2020)*, pp. 1007–1015.