VisionPulse (दिव्य दृष्टि) | ೧ Github

Al-Enhanced Learning and Interaction Assistant for Visually Impaired people at IIIT Delhi

M. Mittal U. Venaik A. Kushwaha L. Kumar T. Singh S. Kabra 2021538 2021570 2021514 2021061 2021569 2021563

Abstract

"VisionPulse" - our revolutionary idea in the educational domain to cater to the needs of visually impaired students at IIIT Delhi by leveraging multiple facets of Computer Vision and Natural Language Processing techniques to build a Large Multi-Model (LMM) that will mitigate the challenges faced by them & bridging the gap in their educational realm by providing them with highly tailored, context-aware assistance system tailored to work in the college campus setting. This will enable them to interact seamlessly with their peers and other objects within the campus environment. Our model will be powered with facial recognition techniques, deciphering complex scenarios in a localized environment and delivering precise, actionable knowledge about the surroundings to turn daily campus navigation and social encounters into a vivid experience. VisionPulse uplifts the barrier of passive participation in classroom environments by leveraging various techniques to extract audio from lecture recordings and to create them into structured, accessible notes, which will further facilitate dynamic, content-specific Q & A sessions using Retrieval Augmented Generation (RAG). The benefits of this technology ensure that no student feels alienated due to specific barriers in their education by enhancing their class participation through information retention techniques. sides this, our project aims to provide the user with more features like realtime obstacle detection, spatial location detection, aiding in reading emails, reading, understanding, and conveying text-based instructions in front of the user, enhancing the feasibility and easing the individual's life. VisionPulse will be able to perform actions based on voice commands and generate audiobased responses as well. Our contribution of VisionPulse to the IIIT Delhi community can act as a gateway in the realm of education for visually impaired students by creating a more inclusive, empowering, and independent educational experience for them. Our project has LLM Project Proposal 2 the potential to become one of the breakthrough contributions in the academic landscape which leverages existing tools and technology to a significant effect, which will help bridge the gaps in education for differently-abled students and create a healthy, inclusive environment where every student can thrive & unlock their full potential to scale great heights & achieve all their dreams.

1 Introduction

Incorporating the emotional dynamics of humans alongside modern methods of computational linguistics and natural language processing is a method which promises significant advancements in how machines understand and interact with humans. Our project addresses this challenge by developing a model capable of pinpointing the exact causes of emotions during conversations. Whether emotions arise from specific statements or are influenced by inherent personal feelings, our model aims to accurately predict these causes and triggers. Our model integrates cutting-

edge NLP techniques like Transformers, attention mechanisms and Convolution Deep Neural networks (CNNs) alongside psychological insights—specifically, the Myers-Briggs Type Indicator (MBTI)—our approach predicts the emotions, emotional causes and triggers based on the conversation. Please refer to background and problem statement sections below for examples of what the model is predicting. Please refer to the problem statement section below for example.

2 Motivation

- 1. Visually impaired people face a lot of challenges doing usual tasks & are often alienated from activities which most of us can do without much fuss. Hence, the primary motivation of our work was to empower the visually impaired people to do everyday tasks at ease & help them overcome the challenges that they face often.
- 2. Inform the user whether they can move forward and provide real-time alerts, such as beeping, when obstacles are detected.
- 3. Since the visually impaired person may not be able to access documents like IIIT Delhi policies, etc., we developed a model that has the knowledge related to IIIT Delhi built within it, which can be accessed and queried by the user at ease.
- 4. Detect the user's current location and provide guidance on routes to reach their desired destination.
- 5. Integrate the student's mailbox by developing an email agent integrated with SMTP to enable users to read and write emails using audio prompts eg: "write a formal mail to admin B.Tech inquiring about the last date of fee payment".
- 6. Allow users to control and navigate the system through voice commands. We aim to tailor each feature of this system to enhance its usability on the IIIT Delhi campus.

3 Problem Statement

Visually impaired students at IIIT Delhi face a significant challenge: the lack of accessible, context-aware educational support tailored to their unique needs. Existing tools often fall short, limiting their ability to seamlessly interact with peers, instructors, and their physical environment. This gap restricts their participation in both academic and social activities.

VisionPulse's Solution:

VisionPulse seeks to address this challenge by harnessing advanced Computer Vision and Natural Language Processing techniques to create a Large Multi-Modal (LMM) system. This innovative solution aims to provide visually impaired students with:

- Real-time obstacle detection
- Location Identification (Visual Question Answering)
- Agent-based system for reading and sending emails
- LLAMA 3 Guard
- General Chat tool
- Student ERP Portal Agent (SQL Agent)
- Agent for retrieving IIIT Delhi Policies

4 Related Work

4.1 VQAsk - A Multimodal Android GPT-based Application

The authors discuss their developed application, VQAsk, available on Android, designed to assist visually impaired users by answering questions about their environment through speech-based interaction and haptic feedback. The app employs MiniGPT-4, a vision-language model, and uses automatic image segmentation for object identification. The system allows users to interact with their surroundings through voice commands and enables visual question answering through images. The dataset used in this study, VisWiz, consists of 31,000 visual questions paired with ten answers as ground truth, making it robust for training and fine-tuning the model. (De Marsico et al., 2024)

4.2 Visual Question Answering for Visually Impaired People (VizWiz-VQA)

This model uses the VizWiz-VQA dataset, which consists of visual questions asked by visually impaired users. This dataset differs

from others in that it is very close to real-world scenarios since the images and answers to questions are taken by the visually impaired people themselves, and the questions are answered based on real-life questions they may have. The study highlights the unique challenges that arise from this authentic use case. Since people face challenges in taking the right pictures, the dataset has various issues like blurred images, which may be encountered in real life. (Chen et al., 2022)

4.3 Wu-Palmer Similarity Score (WUPS Score)

The WUPS score evaluates system-generated answers by using Fuzzy Set theory and Wu-Palmer Similarity to account for semantic fuzziness between classes. It measures the similarity between predicted and actual answers, awarding partial credit when answers are semantically close. This is particularly valuable when exact matches are rare, but approximations are still informative. WUPS penalizes both underestimation and overestimation, ensuring balanced evaluations and providing a more realistic assessment of performance, especially in tasks like visual question answering where ambiguity in data is common. (Malinowski and Fritz, 2014)

4.4 LLM-Assisted VQA Evaluation (LAVE)

Developed LAVE evaluation metric, which uses LLMs like Flan-T5 and GPT to score candidate answers based on their semantic similarity to reference answers instead of using legacy techniques like VQA and soft VQA, outperforming traditional metrics like VQA Accuracy and BERTScore. It was demonstrated that the LAVE metric has higher correlation to human-based evaluation than BERTScore and VQA score. (Mañas et al., 2024)

4.5 Young Persons with Visual Impairment: Challenges of Participation

This paper identified many areas via real-life surveys wherein visually impaired youth face difficulties. Key facets recognized are independent mobility, social isolation, and barriers to accessing information. The study highlighted that visually impaired students often struggle to navigate through educational environments without external support, maintain social connections, and access, as well as understand and interpret written materials with ease. (Salminen and Karhula, 2014)

5 Assistant Features

5.1 Real-time Obstacle Detection

The objective is to develop a model that helps any visually impaired individual detect any obstacle in his/her path while walking, so our model will take an image in the input and classify it as **obstacle** or **Non-obstacle**.

5.1.1 Dataset Collection & Description

Manually collected various images of different locations of IIIT Delhi focusing on obstacles visually impaired individuals may encounter at different body levels.

Head-Level Perspective Images captured at head height to simulate the natural viewpoint, focusing on high-standing obstacles like door frames, signs, branches, etc.

Knee-Level Perspective Images captured at knee height to identify low-lying obstacles such as chairs, steps, small animals, etc.

- Training Dataset: 167 images of IIITD
- Class 0 (Obstacle): 80
- Class 1 (No obstacle): 87
- Testing Dataset: 40 images of IIITD

5.1.2 Model Architecture and Approach

Teachable Machine: CNN-based model

- Input (224x224x3) -> 3 CNN layers -> Output (2 classes)
- The preprocessing steps involve resizing the image to 224x224 dimension and normalizing the pixel values between 0 and 1 after dividing by 255.
- It is compiled with the Adam optimizer, binary cross-entropy loss, and accuracy as the metric.

This is an attempt to use a very rudimentary CNN-based model to give a sanity
check that this problem requires an acute
understanding of the object in the model.

5.1.3 Results

The below table shows the results of the training and testing set:

Metric	Training Set	Testing Set
Accuracy	0.9401	0.8750
Total Inference Time	0.6435	0.1614
Avg Inference Time	0.0039	0.0040
False Negative Rate	0.0805	0.0000
Recall	0.9195	1.0000

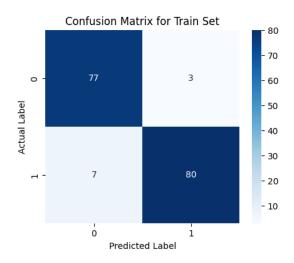


Figure 1: Confusion Matrix (Training)

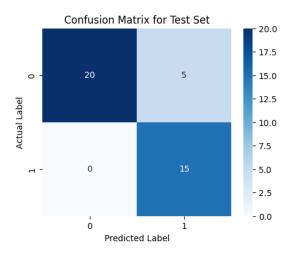


Figure 2: Confusion Matrix (Testing)

5.2 Receive Emails Tool

To enhance communication and accessibility for visually impaired individuals, we integrate an IMAP server connected to their Gmail account. The system will notify the user of new emails, read the latest email aloud, and convert its content into speech for easy consumption by visually impaired students.

- Real-Time Email Notifications:
 Users are notified instantly when new
 emails arrive, ensuring they remain up to
 date without needing to manually check
 their inbox.
- Email-to-Speech Conversion: The latest email is read aloud using text-to-speech, allowing visually impaired users to easily access email content without visual input.
- Query-Based Interaction: The system is designed to process user queries related to email content. The Llama 3.1 orchestrator classifies which tool or functionality to use based on the query, ensuring that the right action is taken (e.g., reading specific sections of emails or answering specific questions based on the email content).

Input:



Llama orchestrator: Selects Appropriate tool according to the input query from all available tools.



Once an email is received, it is read aloud, providing users with easy access to email content without needing to see the screen.

```
Sender-Assah tushanha
Mail Address: asaskY2014iiitdacin
Südject: bip meeting
Hi Utkarsh Can we have our bip meeting today at 5 pm best regards Akash Kushwaha student council 2425 2021314 csai25
```

5.3 Send Emails Tool

The Send Emails Tool integrates an LLAMA model, offering AI-powered functionality to process and send emails efficiently. This tool is designed to streamline email-related queries and automate the process of crafting and sending emails.

- LLAMA Model Integration: The tool employs a LLAMA model for interpreting user queries and classifies them to determine if they are email-related.
- Query Classification: The LLAMA Orchestrator classifies the nature of the query (e.g., "send email") and invokes the appropriate tool accordingly.
- Message Generation: The tool accesses the user's mailbox using SMTP and sends the email to the recipient with the help of the LLAMA model. The LLAMA model incorporated at this step, can both send the exact content written by the sender to the recipient or able to generate a message using a query from the sender (For eg: "write an in-depth Diwali message to abc@gmail.com".

The query from the sender can be made using voice commands, which will be captured, processed & directed to the LLAMA model for target generation.

Input:



Figure 3: User Interaction

The email is sent to the specified recipient through the integrated SMTP service, as confirmed by the tool.



Figure 4: Terminal model output: Mail Generation

5.4 Shared Rolling Memory

Motivation: The Assistant must have a persistent memory of the user queries.

The Orchestrator uses different models and tools that use API calls; thus, there are no internal memories of these queries.

We implemented a shared memory for the orchestrator and down-flow agents and tools to interoperate with the context of previous queries and their output by the orchestrator and downstream tools.

To make a rolling persistent memory of user queries and different tools' answers, Store the last 50 queries and respective answers in textual format. The number of queries taken as 50 is a hyperparameter that can be tweaked based on the user's available storage.

This contextual knowledge base is referenced by the orchestrator and other tools for future queries to keep the AI Assistant in context and give better results.



Figure 5: Shared Rolling memory

5.5 Llama Guard

Motivation: To ensure user safety. We implemented *Llama Guard 2* as a safeguard tool. Its purpose is to prevent harmful content by filtering both user inputs and AI-generated responses, addressing risks such as violence, hate speech, privacy violations, and exploitation.

Model Overview: Llama Guard 2 is an 8B parameter Llama 3-based LLM trained to classify and mitigate harmful content based on the MLCommons taxonomy of hazards, encompassing 11 harm categories:

- 1. Violent Crimes
- 2. Non-Violent Crimes,
- 3. Sex-Related Crimes,
- 4. Child Exploitation,
- 5. Specialized Advice,
- 6. Privacy,
- 7. Intellectual Property,
- 8. Weapons,
- 9. Hate,
- 10. Self-Harm,
- 11. Sexual Content

Approach:

- 1. Query Filtering: User queries are checked by the Gaurd tool and flagged if any potential harm.
- 2. **Agent Processing:** Safe queries are routed to specialized agents for appropriate responses by the orchestrator.
- 3. Response Filtering: Generated responses are re-evaluated and flagged if harmful, ensuring safety before delivery to users.

This dual-layer safeguard effectively addresses harmful queries and generated content.



5.6 General Chat Tool Functionality

The **General Chat Tool**, implemented as one of the tools invoked by the LangChain agent, incorporates a **LLAMA model** in its implementation.

Key Features:

- 1. Fact-based Responses: The LLAMA model is adapted to provide generic, fact-based responses to the user based on the queries fed into the system. It mimics the functionality of a general chat-based mechanism that interacts with the user and provides responses based on any kind of prompt or query.
- 2. Adaptable Query Handling: The tool allows the user to:
 - Request basic overviews of the system.
 - Explore functionalities of the VisionPulse model, such as Real-Time Object Detection.
 - Ask for a wide range of IIIT-Delhibased queries or other diverse topics of user interest.

3. System Integrity Backup: It acts as a fallback mechanism to maintain the integrity of the system. If a specific query aimed at performing a specific task fails, either due to issues in invoking a tool or retrieving a response, the General Chat Tool provides a generic response to the user. This backup mechanism activates after the system reaches the threshold for the maximum number of retry attempts allowed (15 retries, to be precise).

Thus, the General Chat Tool ensures robustness between UI interactions between the system and the user, by providing adaptable responses and serving as a fallback mechanism to maintain system reliability in case of any potential failures.

5.7 Student ERP Portal Agent Functionality

The **Student ERP Portal Agent** is highly useful as it has access to the entire ERP system of the user and can provide a variety of responses based on user queries.

Key Features:

Agent within an Agent Framework:
 The Student ERP Portal Agent contains an SQL Agent within it, designed to provide responses to Natural Language Queries provided by the user.

2. Query Processing and Retrieval:

- The query from the user is processed and sent to the SQL agent.
- The SQL agent formats it into a syntactically correct and coherent SQL query for retrieving the specific information requested, by searching that in the ERP database.
- 3. Response Formatting: The retrieved response, originally structured in SQL format, is:
 - Parsed into a suited format.
 - Propagated back to the main agent.
 - Depicted in the desired format to ensure user readability.

4. Example Use Cases:

- (a) "Tell me my grade in 1st semester."
- (b) "Tell me my grade in NLP course."
- (c) "How many credits have I completed till now?"
- (d) "What were the courses where I got a perfect 10?"
- (e) "What are my 2 worst grades and in which course?"

Visualization: Below is an example query from the user to the Student ERP Portal Agent, where it is invoking the SQL Agent to generate responses based on the query:



Figure 7: Example Use-Case of Student ERP Portal Agent

Thus, the Student ERP Portal Agent improves the user experience by allowing natural language interaction with the ERP system, thereby, ensuring efficient and accurate retrieval of information from the ERP portal. This use case is particularly useful for Visually Impaired Students as they can ask queries using voice-commands and retrieve information in the same format, which helps them to keep tab of important information related to their academics.

5.8 Python Agent Tool

The Python agent tool, referred to as 'python_tool,' integrates VisionPulse access to a wide range of essential and advanced Python libraries. It can perform diverse operations, working on natural language queries provided by the user. This tool can perform tasks like listing files within a directory, generating text files based on user prompts (e.g., creating a text file to track deadlines for a course), and executing complex Python scripts, including arithmetic operations like long-digit multiplications or additions. This



Figure 8: Python Agent Tool

tool can be connected both to the laptop and to the server.

It's applications are extensive. It can be used to generate text files or PDFs, execute system-level updates, solve mathematical problems, etc. It eliminates the need to visually navigate and manage files and other tasks, and it bridges the gap in accessibility.

References

[Chen et al.2022] Chongyan Chen, Samreen Anjum, and Danna Gurari. 2022. Grounding answers for visual questions asked by visually impaired people. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19098–19107.

[De Marsico et al.2024] Maria De Marsico, Chiara Giacanelli, Clizia Giorgia Manganaro, Alessio Palma, and Davide Santoro. 2024. Vqask: a multimodal android gpt-based application to help blind users visualize pictures. In *Proceedings of the 2024 International Conference on Advanced Visual Interfaces*, pages 1–5.

[Malinowski and Fritz2014] Mateusz Malinowski and Mario Fritz. 2014. A multi-world approach to question answering about real-world scenes based on uncertain input. Advances in neural information processing systems, 27.

[Mañas et al.2024] Oscar Mañas, Benno Krojer, and Aishwarya Agrawal. 2024. Improving automatic vqa evaluation using large language models

[Salminen and Karhula2014] Anna-Liisa Salminen and Maarit E Karhula. 2014. Young persons with visual impairment: Challenges of participation. Scandinavian Journal of Occupational Therapy, 2014:1–10.