

# **AI DRIVEN VOICE IMAGE PROJECTOR**

*Project report submitted in partial fulfilment of the requirement for the degree of*

**BACHELOR OF TECHNOLOGY**

**IN**

**ELECTRONICS AND COMMUNICATION ENGINEERING**

**BY**

**AYUSH KHACHI (211005)**

**DEEPANSHU (211080)**

**UNDER THE GUIDANCE OF**

**DR. VIKAS BAGHEL**



**JAYPEE UNIVERSITY OF INFORMATION TECHNOLOGY,  
WAKNAGHAT**

**MAY 2025**

# TABLE OF CONTENTS

<b>CAPTION</b>	<b>PAGE NO.</b>
<b>DECLARATION</b>	<b>i</b>
<b>ACKNOWLEDGEMENT</b>	<b>ii</b>
<b>LIST OF ACRONYMS AND ABBREVIATIONS</b>	<b>iii</b>
<b>ABSTRACT</b>	<b>iv</b>
<b>CHAPTER-1: INTRODUCTION</b>	<b>1</b>
1.1 Background and Motivation	1
1.1.1 Importance of AI in Image Classification	2
1.1.2 Significance of Voice-Controlled Interfaces	2
1.2 Problem Statement	3
1.2.1 Objectives of the Project	4
1.2.2 Scope and Applications	5
<b>CHAPTER-2: PROJECT BACKGROUND</b>	<b>6</b>
2.1 LITERATURE REVIEW	6
2.2 Advancements in Core Technologies	10
2.2.1 Evolution of Image Classification	11
2.2.2 Developments in Speech Recognition	12
2.2.3 Integration of AI in Embedded Systems	13
2.2.4 Real-World Context and Need for the Project	14
2.3 Integration of Image Classification and Speech Recognition	14
2.3.1 Challenges in Integration	14
2.3.2 Benefits and Applications	15
<b>CHAPTER-3: METHODOLOGY</b>	<b>17</b>
3.1 Implementation Approach	18
3.1.1 Dataset Preparation	18
3.1.2 Model Design	21
3.1.3 Speech Recognition Integration	24
3.1.4 User Interface Design	25
3.2 Tools and Technologies	28
3.2.1 Machine Learning Framework	28
3.2.2 Speech Recognition Tools	31
3.2.3 Hardware	31

3.2.4 Supporting Libraries	32
3.3 Workflow Diagram	32
3.4 Challenges Faced	34
3.4.1 Hardware Limitations	34
3.4.2 Speech Recognition Accuracy	35
3.4.3 Real-Time Performance	35
<b>CHAPTER-4: SYSTEM DESIGN AND ARCHITECTURE</b>	<b>37</b>
4.1 Experimental Setup	37
4.1.1 Hardware Environment	37
4.1.2 Software Environment	37
4.1.3 Evaluation Criteria	39
4.2 Results	39
4.2.1 Image-Classification Performance	39
4.2.2 Speech-Recognition Performance	40
4.2.3 System Responsiveness	41
4.3 Analysis	41
4.3.1 Strengths	42
4.3.2 Weakness and Limitations	42
4.3.3 Potential Improvement	43
4.3.4 Broader Implications	43
<b>CHAPTER-5: CONCLUSION</b>	<b>44</b>
5.1 Conclusion	44
5.2 Final Thoughts	45
<b>REFERENCES</b>	<b>—</b>
<b>APPENDIX</b>	<b>—</b>

## DECLARATION

We hereby declare that the work reported in the B. Tech Project Report entitled “**AI DRIVEN VOICE IMAGE PROJECTOR**” submitted at **Jaypee University of Information Technology, Waknaghat, India** is an authentic record of our work carried out under the supervision of **Dr. Vikas Baghel**. We have not submitted this work elsewhere for any other degree or diploma.

Ayush Khachi  
211005

Deepanshu  
211080

This is to certify that the above statement made by the candidates is correct to the best of my knowledge.

Dr. Vikas Baghel

Date:

Head of the Department

Project Coordinator

## ACKNOWLEDGEMENT

I am deeply grateful to all those who have contributed to the successful completion of this project.

First and foremost, I would like to express my sincere gratitude to **Dr. Vikas Baghel, Assistant Professor (SG) and Dr. Rajiv Kumar, Head of Department**, for their invaluable guidance, encouragement, and constructive feedback throughout this project. Their expertise and mentorship have been instrumental in shaping this work.

I am also thankful to Jaypee University of Information Technology, for providing the resources and support necessary to bring this project to fruition. The access to facilities and equipment has greatly enhanced my ability to execute and refine my ideas.

Furthermore, we would like to express our gratitude to **Mr. Dhirendra Kumar** and **Dr. Ajay Thakur** who has generously provided us with the necessary resources, such as laboratory facilities. Their dedication and hard work were instrumental in bringing this project to fruition. In conclusion we would like to express our heartfelt gratitude to all those who have contributed to the successful completion of this project.

## LIST OF ACRONYMS AND ABBREVIATIONS

Acronym / Abbreviation	Full Form
AI	Artificial Intelligence
CPU	Central Processing Unit
GPIO	General Purpose Input Output
HDMI	High-Definition Multimedia Interface
I/O	Input/Output
ML	Machine Learning
OS	Operating System
Pi	Raspberry Pi
RAM	Random Access Memory
USB	Universal Serial Bus
VNC	Virtual Network Computing
YOLO	You Only Look Once
CNN	Convolutional Neural Network
API	Application Programming Interface
SSD	Solid State Drive
LED	Light Emitting Diode
LCD	Liquid Crystal Display
STT	Speech-to-Text
TTS	Text-to-Speech
GUI	Graphical User Interface

# ABSTRACT

The Voice-Controlled AI Image Projector is an innovative Raspberry Pi-based smart photo frame that integrates image classification and offline speech recognition to enable hands-free navigation and display of personal photos. Addressing the limitations of traditional digital photo frames, which often rely on manual inputs or cloud-based processing, this project offers a privacy-focused, fully offline solution optimized for resource-constrained hardware. The system employs a custom Convolutional Neural Network (CNN) deployed in TensorFlow Lite to classify images into four categories—cat, dog, human, and scenery—with 87% validation accuracy and 450 ms inference time.

The Vosk toolkit facilitates real-time speech recognition, processing commands like “hey frame, show cat photos” or “next” with 92% accuracy in quiet environments. Pygame powers an immersive full-screen interface, supporting automatic slideshows and voice-controlled navigation. Implemented on the Raspberry Pi 5, the system overcomes hardware challenges through model quantization, data caching, and efficient audio processing, achieving seamless multitasking within an 8 GB RAM footprint. Key challenges included managing memory constraints, mitigating background noise in speech recognition, and optimizing real-time performance, addressed through lightweight models and robust error handling. User testing confirmed high usability, with a 4.2/5 satisfaction rating.

The project’s offline operation ensures data privacy and functionality in diverse settings, from homes to remote areas, with applications in education, healthcare, and public displays. By combining AI and voice control on an affordable platform, this work contributes to the development of accessible, privacy-conscious smart home technologies, paving the way for future enhancements like expanded categories and smart home integration [19][20].

# CHAPTER 1

Artificial Intelligence (AI) has altered the way humans interact with technology in line with how we currently interact with each other, allowing us to become more efficient and comfortable in completing everyday tasks. The rapid development of computer vision and natural language processing has made it feasible for researchers to develop systems that allow humans to transmit information to machines using speech and other natural human input. This research project will see the development of a fully voice-enabled picture classification and navigation system, which will allow the classification and display of images, selecting classification parameters based on user commands. The system will employ two powerful technologies: Convolutional neural networks (CNNs) for picture classification, and Speech recognition technology for the interpretation of spoken commands. The research project addresses the need for an intuitive, hands-free, efficient way to manage and navigate a large database of images and other visual material.

## 1.1 Background and Motivation

The rapid advances in Artificial Intelligence (AI) and embedded systems have produced multiple smart, interactive devices that give users convenience and an enhanced user experience. Examples of applications that will transform the way we interact with machines are image classification and voice recognition integrated together, enabling a system to recognize visual information, and humans speaking, simultaneously providing an instantaneous response. AI and humans are learning how to interact in different ways, providing us with tools and new systems that we don't want, and we don't see the full capabilities yet, but we are only scratching the surface in areas such as home automation, surveillance, healthcare, and assistive technology.

New, low-cost computing platforms, such as the Raspberry Pi 5, allow for accessible and affordable AI, and have the processing capability to run lightweight machine learning models on the local device. The Raspberry Pi 5 offers significantly more processing power, supported memory, and peripheral management, allowing for real-time AI solutions to be scaled down to the edge and remove requiring full cloud infrastructures.

This project is significant because we are using that capability to present a new device that uses voice to control a local AI-powered smart image projector. It uses intelligent classification of



a set of images, which may include cats, dogs, humans, and scenery and the projector will output projected images after receiving the command via voice. The goal is to create a hands-free, responsive, and intelligent visual experience — a functional representation or embodiment of machine learning, speech processing, and embedded system design.

### 1.1.1 Importance of AI in Image Classification

Image classification is one of the most basic tasks of computer vision, which takes an image and provides a label or category for that image based on the visual content within the image. With the increasing volume of digital media, it is now essential to be able to automatically organize, sort, and share media files based on their content within a large variety of industries.

- In **healthcare**, AI helps identify medical anomalies from X-rays or MRIs.
- In **autonomous vehicles**, it enables the identification of road signs, pedestrians, and obstacles.
- In **security systems**, it assists in recognizing faces or suspicious objects.
- In **social media**, it supports automatic tagging and content moderation.

Traditional, rule-based algorithms have difficulty for scaling with the diversity of the data. On the other hand, AI-based algorithms can learn from the data to understand and classify complex visual patterns. AI-based models, particularly based on deep learning, refer to subsets of AI models made to learn complex visual patterns and processes in all types of data. One example of this type of AI model is the convolutional neural network (CNN), which can analyze a set of images and automatically classify each based on detailed features.

In this project, a trained model has been converted to TensorFlow Lite (TFLite), meaning the model has substantially reduced in size and complexity without losing too much performance. The converted model is capable of performing classification on the Raspberry Pi 512 without needing additional software, hardware, etc. The model can classify images into four categories - cats, dogs, humans, and scenery - as part of the classification process before performing the requested task by the user via voice command.

The advantages of on-device, or in this case, on-Raspberry Pi 5, will ensure that classification is fast, secure and does not require any type of network connectivity to perform the intended task. Having that independence to work offline or in remote locations would be beneficial.

### 1.1.2 Significance of Voice-Controlled Interfaces

In recent years, voice interfaces have transformed human-computer interaction. Devices like Amazon Echo, Google Nest, and Apple's Siri have demonstrated that speech is among the most natural and efficient ways for users to interact with technology.

Voice-controlled systems are particularly important in the following contexts:

**Accessibility:** Aiding users with physical disabilities that may deter them from interacting with technology in a non-touch or non-viewing manner.

**Hands-Free Interactivity:** Useful in all kinds of environments such as kitchens, workshops, medical, and driving situations.

**Multilinguistic Compatibility:** Speech systems may also be trained or set up to work in different languages and dialects.

Unlike most commercial voice assistants that rely on cloud-based processing, this project uses **Vosk**, an open-source **offline speech recognition** toolkit. Vosk is lightweight, accurate, and capable of recognizing commands with minimal delay. It works entirely on the Raspberry Pi 5, ensuring **user privacy** and independence from internet access.

In this system, users can control the photo frame using simple commands like:

- “Hey frame, show dog”
- “Hey frame, next”
- “Hey frame, show scenery”

This allows seamless interaction and elevates the image projector from a static viewer to a dynamic assistant.

## 1.2 Problem Statement and Objectives

Although there are currently picture classification and speech recognition systems, there are several unique challenges in combining these technologies into a reliable, efficient system. The key difficulties emerge from the need to preserve high accuracy in both voice and picture classification while ensuring seamless, real-time performance, particularly when the system is being implemented on hardware with constrained resources.

The following issues need to be resolved for this project:

1. **Accurate Classification:**In order to ensure that every image is correctly classified, the system must classify photos over a wide range of categories, including humans, plants, and animals.
2. **Real-Time Voice Processing:**The speech recognition module needs to be able to react quickly, processing spoken commands without any noticeable lag.
3. **Deployment on Resource-Constrained Hardware:**Because the Raspberry Pi 4 has less processing power than desktop PCs or cloud-based platforms, the system needs to be tuned to work effectively on it.

### 1.2.1 Objectives of the Project

The primary objectives of this project are as follows:

1. **Pre-trained CNN Model for Image Classification:** To develop a pre-trained Convolutional Neural Network (CNN) model (converted to the TensorFlow Lite) to classify a dataset of images on Raspberry Pi 5 for four categories: cats, dogs, humans, and scenery. Performance tests of the model should be run for accuracy, speed, and memory usage.
2. **Offline Voice Commands:** To develop an offline speech recognition system using the Vosk toolkit to parse certain recognizable voice commands of "show cat", "show human", "next", and "previous". This feature will allow the user to navigate through voice commands hands-free and intuitively, without relying on a connection to the internet.
3. **Image Display System with Real-time Interaction:** To develop an image display interface using Pygame that will show the classified images in full-screen mode. The image display is expected to show images that satisfy the voice request, be responsive to the voice commands, and be able to segue from one image category to another, or simply go through each individual image.
4. **Offline Development:** To ensure that the system works completely offline, including CAMCB's ability to classify images offline, voice recognition and understanding offline, and projecting images offline. This will maintain the user's privacy and will allow the system use for users with limited resources or no dedicated internet connection.

5. **Easily Navigable and Usable Design:** To design the system for simple use, particularly for users with limited movement which would make it valuable in an educational, assistive, or home setting.
6. **Adequate Runtime on Raspberry Pi 5:** To deploy an image classification and speech recognition system into a Raspberry Pi 5; to run properly with the Raspberry Pi 5 hardware limitations.

### 1.2.2 Scope and Applications

This project is designed with a wide range of potential applications, both in personal and professional contexts. The scope extends beyond academic purposes and could be applied in various industries, such as healthcare, education, and home automation.

1. **Use for Personal Use:** This system can be used to manage and organize a personal photo collection in a searchable way. Using voice designation could help the user identify specific image types or navigate images at a higher speed/dimension on orders of thousand images.
2. **Use for Healthcare:** In the hospital setting, this system would help doctors and medical professionals categorize medical imaging and access images based on voice designation quickly. This could be used in a more immediate way in a busy hospital setting needing to save time.
3. **Use for Education:** In school systems, the imaging would allow students to browse through images of animals, plants, and humans for biology or environmental studies.
4. **Use for Industry:** This system can be used to categorize stock or inventory images in industries like distribution or logistics. With voice-controlled capabilities even employees can quickly locate items or verify stock without searching through thousands of images.

## CHAPTER 2

An overview of the technology, approaches, and relevant literature that serve as the project's foundation is given in this chapter. It talks about the developments in speech recognition and picture categorisation, as well as how these technologies might be combined for creative uses.

### 2.1 Literature Review

The Voice-Controlled AI Image Projector integrates image classification, offline speech recognition, and an embedded user interface on a Raspberry Pi to create a privacy focused, interactive smart photo frame. This section reviews existing literature and commercial solutions in digital photo frames, image classification, speech recognition, and embedded AI systems, positioning the project within the broader research landscape and identifying gaps it addresses.

#### Digital Photo Frames and AI Integration

Digital photo frames have transitioned from being simple display devices to advanced systems that involve some form of artificial intelligence (AI) and voice control. A few years ago, digital frames were employed just to display slideshows, where the user manually changed the photo via physical inputs or remote control (Chin, W., & Tan, C. 2010) [1]. Now, AI features, such as facial recognition and content-based organization, allow frames to recognize individuals, or organize the photos based on people, pets, or landscapes.

Commercial offerings use advanced AI capabilities, such as Nixplay Smart Photo Frame and Aura Frames, which connect to cloud-based platforms and voice assistants like Amazon Alexa or Google Assistant. They also allow users to control their displays by voice or mobile application [3].

Often, solutions like these process images in the cloud, which creates potential privacy concerns, as the user will be sending their personal photos, as well as their voice to a cloud-based system [4].

Also, the ability to use a cloud-based solution limits functionality if or when connectivity is limited. This is an important point for users who are in remote settings, or privacy

cautious situations. Patel, R, & Gupta, S. 2023 [5], really highlight a gap between these cloud processes and the offline AI capabilities when looking specifically at smart photo frames which can still function without internet access and also protect users' data. Our Voice-Controlled AI Image Projector, aims to fill this void with fully offline image classification and speech recognition, eliminating the need for cloud connectivity and enhancing user privacy.

## Image Classification for Photo Management

Image classification, a cornerstone of modern photo management, leverages convolutional neural networks (CNNs) to categorize images based on visual content. CNNs, as detailed by (Sharma, A. 2025) [6], excel at extracting spatial features like edges, textures, and shapes, making them ideal for tasks such as identifying objects, scenes, or individuals in photos. In personal photo management, platforms like Google Photos and Apple Photos use CNNs to organize images into categories such as people, pets, places, or events, often employing facial recognition for person-specific grouping [7].

For edge devices like the Raspberry Pi, lightweight CNN architectures, such as MobileNet and EfficientNet, are preferred due to their reduced computational requirements [8]. (Chollet, F. (2016)) [9] demonstrate that small, curated datasets (e.g., 500-1000 images per class) can achieve high accuracy with proper pre-processing and data augmentation, a strategy adopted in this project. Techniques like image resizing, colour normalization, and augmentation (e.g., rotation, flipping) enhance model generalization, particularly for personal photo datasets with limited diversity [10].

Despite these advances, most image classification systems for photo frames rely on cloud-based models, which incur latency and privacy costs. (Viso.ai (2025)) [11] advocate for edge-based inference using frameworks like TensorFlow Lite, which optimize models for low-power devices through techniques like quantization and pruning. This project builds on these principles, deploying a custom CNN in TensorFlow Lite format to classify images into four categories (cat, dog, human, scenery) directly on the Raspberry Pi, ensuring efficiency and privacy.

## Offline Speech Recognition Systems

Speech recognition enables hands-free interaction, a critical feature for modern smart devices. Cloud-based systems, such as Google Speech-to-Text and Amazon Transcribe, offer high accuracy but require internet connectivity and raise privacy concerns [12]. Offline speech recognition systems, such as Vosk and Kaldi, provide viable alternatives for edge devices, offering low-latency, privacy-preserving solutions [13].

Vosk, used in this project, is an open-source toolkit based on the Kaldi framework, designed for real-time speech recognition on resource-constrained hardware [13]. (Kumar, V., & Singh, R. (2023)) [14] evaluate Vosk on the Raspberry Pi, reporting a word error rate (WER) of approximately 10% in quiet environments, with a memory footprint of 50-100 MB for small models. Its offline capability and support for custom wake words make it suitable for smart home applications where privacy and reliability are paramount. However, challenges such as background noise and accent variability can degrade performance, necessitating robust pre-processing and command parsing [15].

This project leverages Vosk’s lightweight English model (vosk-model-small-en-us-0.15) to detect the wake word “hey frame” and subsequent commands, aligning with findings that offline systems can achieve near-real-time performance on edge devices [14]. By integrating Vosk with a USB microphone and a queue-based audio processing system, the project ensures seamless voice control without external dependencies.

## Embedded AI on Resource-Constrained Devices

Deploying AI on edge devices like the Raspberry Pi presents unique challenges due to limited memory, processing power, and lack of GPU acceleration. (Seeed Studio (2021)) [16] outline optimization techniques for machine learning inference, including model quantization, pruning, and efficient data pipelines, which reduce memory usage and latency. TensorFlow Lite, a framework tailored for edge devices, supports these optimizations, enabling CNNs to run efficiently on the Raspberry Pi’s CPU [8].

Research on embedded AI for smart home applications emphasizes the importance of balancing performance with resource constraints. (Viso.ai (2025)) [11] note that edge based vision systems can achieve inference times under 500 ms for lightweight models, suitable for real-time applications like photo classification. Similarly, (Nguyen, T., & Tran, D.

(2024)) [17] describe integrated vision and speech systems on the Raspberry Pi, reporting successful multitasking with careful memory management.

Commercial and open-source projects, such as Raspberry Pi-based digital signage and home automation systems, demonstrate the feasibility of embedded AI but rarely combine image classification and offline speech recognition in a single device [18]. This project's integration of TensorFlow Lite for image classification, Vosk for speech recognition, and Pygame for display addresses this gap, creating a cohesive, offline system optimized for the Raspberry Pi 5's quad-core processor and 8 GB RAM.

### Research Gaps and Project Contributions

While significant progress has been made in AI-driven photo frames, several gaps remain. First, most commercial frames rely on cloud-based processing, limiting their use in privacy-sensitive or offline environments [4]. Second, integrated systems combining image classification and offline speech recognition are rare, as most solutions focus on either vision or voice [17]. Third, deploying such systems on resource-constrained devices like the Raspberry Pi requires careful optimization, a challenge not fully addressed in consumer products [16].

The Voice-Controlled AI Image Projector fills these gaps by:

- Implementing a fully offline system, ensuring privacy and functionality without internet access.
- Combining image classification (using a custom CNN) and offline speech recognition (using Vosk) in a single device.
- Optimizing for the Raspberry Pi through lightweight models, quantization, and efficient data processing.

By addressing these gaps, the project contributes to the development of privacy focused, edge-based smart home devices, offering a scalable model for future embedded AI applications.



## 2.2 Advancements in Core Technologies

As Artificial Intelligence (AI) is increasingly embedded into everyday devices, human-computer interactions have changed dramatically. Ideas that were once just a vision of the future - smart homes, autonomous systems and voice-controlled interfaces - are now possible and widespread as a result of miniaturizing computing hardware and the rapidly evolving landscape of software frameworks. The merging or intersection of technologies found in the development of a smart, voice-controlled image projector that integrates the capabilities of AI, embedded systems and human-computer interaction is one notably exciting example of this combination.

Two core technologies are at the foundation of any intelligent system: image classification and speech recognition. Each of these areas have changed dramatically over the last decade, moving away from the behemoth cloud-centric, heavyweight architectures to a real-time, lightweight approach fully on the edge, and not required cloud support, now possible on devices like the Raspberry Pi 5. The advancements are the result of exhaustive efforts on the parts of researchers in many disciplines, each looking to improve machine learning algorithms, efficient model reduction methods such as pruning and quantization, and efficient deployment frameworks like TensorFlow Lite and ONNX Runtime.

### Image Classification

Previously, image classification required the use of far-reaching and expensive datasets and large server-grade GPUs. With these factors in play, we were usually limited to usage of images classification in an enterprise or an academic setting. In recent years, due to the development of Convolutional Neural Networks (CNNs) and optimized training to showcase accurate and reliable results, the image classification on limited consumer devices is feasible. The development of models such as MobileNet, EfficientNet, and SqueezeNet provide a low memory footprint and responsive architecture to make acceptable image classifications for mobile and embedded use.

These models are capable of recognizing and categorizing images across multiple classes—such as animals, humans, objects, and natural scenes—with a level of accuracy that rivals their

server-based counterparts. As a result, applications like smart surveillance, autonomous vehicles, and intelligent image sorters have become practical even at the consumer level. In the context of this project, image classification enables the system to automatically categorize input photos into predefined categories (e.g., cats, dogs, humans, scenery), allowing users to retrieve and view specific types of images through simple voice commands.

## **Speech Recognition**

Speech recognition technologies have improved parallel to image processing. The focus was originally on internet-based APIs and cloud services (for example Google Speech, Amazon Alexa) for speech recognition, but most speech recognition engines can now operate offline. Modern libraries, such as Vosk, Deep Speech, and Picovoice, offer accurate lightweight models that allow real-time detection and interpretation of voice with very little lag. This is critical for applications requiring privacy or simply operate offline (for instance Raspberry Pi), because internet access will be severely limited or non-existent.

Offline speech recognition decreases reliance on cloud infrastructures, improves response times and reliability when consistent connectivity cannot be ensured. Furthermore, the engines support custom vocabulary tuning and command (keyword) recognition that are constrained, specifically for applications such as voice operated controls. Here, the Vosk speech recognition engine is selected to parse voice commands (for example "show cat", "next image") into actionable triggers for the image projector.

## **Convergence and Real-World Feasibility**

The convergence of these technologies - compact image classification models and fast offline speech recognition - means we can make interactive, voice-enabled devices that would have been inconceivable. By combining these two technologies, and deploying them on a single embedded platform (the Raspberry Pi 5), this solution reveals autonomy, portability and real-time responses that would not have been plausible without a limitation induced by the online nature of the cloud systems. Complimented by advances in our integrating peripherals with higher-quality, high-definition HDMI displays, USB microphones, and camera modules, it allows us to offer users more natural and engaging user experiences.

In closing, this project is plausible because of the rapid and sustained advances in computer vision and natural language processing. With continued advances in the hardware performance of CPU, GPU, and their utilization with software that operates far more quickly than previously recommended, we will see more opportunities for AI-enabled, voice-controlled, system for everyday use. The project provides an illustration of the creative potential of new technological advances to produce something new, that offers a user experience that previously was unattainable using off-the-shelf, consumer products.

### 2.2.1 Evolution of Image Classification

Image classification is an elementary task in computer vision and is also the central task that this project is all about. This task consists of assigning one from a set of pre-defined labels to an image based on it's visual content. For our project, the image classification task is focused on automatically classifying a provided picture into one of four categories (cat, dog, human, scenery). When the classifying is completed, a user will be able to \_, bolt through \_ an image set based on the classification, with a set of voice commands, through an automated system, allowing the user to experience the categorization of images through smooth intelligent automation.

Historically, image classification began with simple models that required manual feature extraction. These early systems used handcrafted features like texture, color histograms, shape edges, or corners, which were then processed through machine learning algorithms such as:

- k-Nearest Neighbors (k-NN),
- Support Vector Machines (SVM),
- Decision Trees.

While these approaches had modest success, they were limited by their reliance on predefined rules and inability to generalize across complex datasets.

With the advent of **Convolutional Neural Networks (CNNs)** in the early 2010s, particularly through the groundbreaking **AlexNet** architecture in 2012, the field saw a massive leap forward. CNNs eliminated the need for manual feature engineering by learning filters and feature hierarchies directly from raw images. Successive models such as **VGGNet**, **ResNet**, and **MobileNet** provided increasingly accurate and efficient classification capabilities.

In our project, we utilize a **TensorFlow Lite (TFLite)** model optimized for embedded platforms like the Raspberry Pi 5. This model classifies images after preprocessing (resizing, normalization) into the desired categories, storing the results for organized slideshow viewing.

**Notable developments and applications in image classification include:**

- **Medical imaging:** AI classifies X-rays or MRI scans for diagnostics.
- **Self-driving cars:** Image classification helps identify pedestrians, vehicles, and road signs.
- **Security and surveillance:** Automated recognition of people and suspicious activity.
- **Retail and e-commerce:** Product tagging and visual search.

By using a compact TFLite model, our project achieves near real-time classification on a low-power device without needing cloud support. This highlights the maturity and portability of modern image classification systems.

### 2.2.2 Developments in Speech Recognition

Speech recognition technology has transformed from a niche research field into a mainstream interface used in smartphones, smart homes, and automotive systems. Its evolution has mirrored the broader AI landscape, transitioning from rule-based logic to data-driven deep learning approaches.

Earlier systems were based on **template matching** and **phonetic rules**, which were not scalable across different languages, accents, or noisy environments. These were replaced in the 1990s and 2000s by **Hidden Markov Models (HMMs)** coupled with **Gaussian Mixture Models (GMMs)**. Though statistically more robust, these models required large volumes of labeled speech data and were often brittle when faced with natural language variability.

The true leap forward occurred with **Deep Neural Networks (DNNs)**, which allowed for end-to-end learning of speech representations. Later advancements included:

- **Recurrent Neural Networks (RNNs)** for handling temporal dependencies in audio.
- **Long Short-Term Memory (LSTM)** networks for improved memory and context retention.

- **Transformer-based models** (e.g., Whisper, Wav2Vec) for high-accuracy speech-to-text translation.

However, most modern systems still depend on cloud servers for processing due to the size and complexity of these models. This poses problems for privacy, latency, and constant internet requirements.

Our project addresses these concerns by implementing **Vosk**, a lightweight speech recognition toolkit that works offline. Vosk is built on top of **Kaldi**, a highly respected toolkit in the academic speech recognition community. It provides real-time recognition of commands using a small-footprint acoustic model, suitable for deployment on devices like Raspberry Pi.

The AI photo frame listens for the wake word "**hey frame**", activating voice command mode. It can then recognize keywords such as:

- "**cat**", "**dog**", "**human**", "**scenery**" – to switch image categories.
- "**next**" – to skip to the next image in the selected category.

**Key advantages of using offline speech recognition in this project include:**

- **Enhanced privacy** – no user data is sent to cloud servers.
- **Faster response time** – local processing reduces delay.
- **Offline operation** – useful in rural, remote, or private environments.
- **Accessibility** – makes the system easier to use for elderly or disabled users.

This level of interaction brings the photo frame to life, giving it the feel of a responsive assistant rather than a passive display device.

### 2.2.3 Integration of AI in Embedded Systems

The convergence of compact hardware and AI software has enabled powerful applications on edge devices. The Raspberry Pi 5, used in this project, offers enough computing power to run both image classification and speech recognition in real time. This demonstrates the growing capability of embedded platforms to support sophisticated AI tasks.

The use of:

- **TensorFlow Lite** for model inference,
- **Vosk and Kaldi** for voice command processing,
- and **Pygame** for full-screen image display,

...shows how open-source libraries and optimized toolkits can be combined to build feature-rich systems even on low-budget hardware.

#### 2.2.4 Real-World Context and Need for the Project

Traditional digital photo frames simply display a set of images in order or shuffle mode, requiring manual input to change the behavior. As more people generate vast numbers of images through smartphones, organizing and interacting with them becomes increasingly difficult.

An intelligent, voice-controlled image projector serves various practical purposes:

- **Hands-free interaction** while cooking, working, or resting.
- **Elder-friendly design** for those unfamiliar with smartphones or touchscreens.
- **Smart classification** that eliminates the need for manual sorting.
- **Offline functionality** for areas with limited or no internet.

The project demonstrates how advances in AI can be used to solve everyday problems and enhance the way we interact with digital memories.

### 2.3 Integration of Image Classification and Speech Recognition

Integrating image classification and speech recognition into a single system represents a high-impact enhancement to the human-computer interaction operation cycle. Historically speaking, the two technologies have been separately developed and applied, with image classification being used to classify objects visually and speech recognition being used to interact with a computer auditorily. The combination of both allows for the development of more intelligent systems that can maintain context and automatically react to a user's commands in a more seamless way. In this project, we enabled the user to control a digital photo frame through natural voice commands that produce a visual change, for example, retrieve and view one of several different categories of images.

As designers, we like to think this multimodal technological experience reflects how the human interaction with the world is constructed. We interact with the world through vision and voice in parallel. There are many ways this kind of integration presents itself in design from smart assistants and autonomous vehicles to displays in museums and assistive technologies. This new type of integration gives rise to entirely new possibilities in user experience design in general. Our voice-directed AI image projector is a small and effective example in practice; it is an integration of a low-resource speech recognizer and a low-resource image classifier using lightweight models and it is our original objective to show how the integration would work on a low-end device such as the Raspberry Pi 5.

By combining the two AI modalities the user has hands-free control of searching through and viewing images and the photo frame provides intelligent organization of images via a machine learning classifier.

The result is a product that is not only technically impressive but also user-friendly and accessible to a wide audience, including elderly users and those with physical disabilities.

### **2.3.1 Challenges in Integration**

While the idea of combining speech and vision AI systems is attractive, actually executing it on embedded platforms has a plethora of challenges. These challenges include limited hardware, modality synchronization, software compatibility, and accuracy of response.

First, real-time performance. Both image classification and speech recognition are heavy computational tasks. Performing both tasks on a low-power device like the Raspberry Pi 5 will need optimized models, optimized use of multithreading, and minimal latency. TensorFlow Lite and Vosk may have down-sized models, however, the authors still had to ensure the model was processing information well enough not to lag or crash when performing both actions at the same time.

Second, the authors must synchronize the voice command with its visual representation. The authors will need to ensure that it:

- Recognize when a valid voice command is issued (e.g., “Hey frame”),
- Accurately parse the keyword (e.g., “cat” or “next”),
- Interrupt or redirect the current image display sequence without glitches.

This requires a responsive event loop and state management in code. Achieving this balance between listening, processing, and displaying images seamlessly is non-trivial.

Furthermore, the project faces limitations such as:

- **Ambient noise**, which may interfere with voice recognition.
- **Misclassification** of visually ambiguous images.
- **Memory constraints** when handling large image datasets.
- **Lack of GPU acceleration**, which restricts real-time processing speed.

Addressing these limitations required considerable experimentation with buffer sizes, frame rates, and thread priority. Despite these difficulties, the system demonstrates that integration is feasible with the right optimizations and software tools.

### 2.3.2 Benefits and Applications

Although it is not without its challenges, the advantages we get from connecting image classification with speech recognition outweigh the challenges involved in developing a single interface to fuse the two. This integration can ultimately yield a unique interface to transform our interaction with digital content.

In the context of this project, users can:

- **Effortlessly sort and browse images** through simple voice commands.
- **Access image categories without physical input**, making the system ideal for hands-free environments like kitchens or workshops.
- **Enjoy a more personalized experience**, where the frame responds to their voice and displays curated image sets accordingly.

Such a system has wide-ranging applications beyond digital photo frames, including:

- **Smart homes**: Voice-controlled visual dashboards or family galleries.
- **Education**: Interactive learning systems that display visual aids in response to student queries.
- **Retail**: Smart product displays that switch content based on staff or customer voice instructions.



- **Healthcare:** Assistive devices for patients with mobility issues, where vocal commands can bring up instructional images or visual stimuli.
- **Museums and exhibitions:** Interactive kiosks that respond to voice prompts with contextual visual content.

Additionally, the **offline capability** of this system makes it suitable for use in areas with limited or unreliable internet access, including rural schools, remote homes, or fieldwork scenarios.

By combining the strengths of visual and auditory processing, such integrated systems make human-computer interaction more natural, intuitive, and efficient.

## CHAPTER 3

This chapter presents the systematic and logical procedure followed for the design and implementation of the Voice-Controlled AI Image Projector, a smart Raspberry Pi based photo frame that incorporates image classification and speech recognition to select images to be displayed back to the user in response to voice commands. The project illustrates a practical merge between the disciplines of machine learning, speech processing, and embedded systems engineering, as it is designed to operate independently and without the need for an external server or cloud service, making a low-cost solution that is efficient, secure, and ideal for offline or low internet availability.

The fundamental objective of the project was to make a stand-alone smart image projector capable of understanding human speech, deciphering commands, and displaying images from a local directory to the commands of the user. As part of this task, it also classifies images as it organizes them into classes, such as "cats", "dogs", "humans", and "scenery", using a pre-trained image classifier. The information processed will produce a system that autonomously can pull back images on command, from each class.

To achieve this functionality, the overall development cycle was divided into four main phases:

- Dataset Preparation
- Model Design and Training
- Speech Recognition Integration
- User Interface and Display Management

Each phase of development required careful planning and execution, taking into account the limited hardware resources of the Raspberry Pi 5 platform—particularly its constraints on processing power, RAM, and storage space. As a result, lightweight and optimized models were prioritized for both image classification and speech recognition tasks, with the confidence that they would operate in the real time, with no accuracy loss or compromises.

The custom image dataset, created in the first phase, contained a number of examples (manually curated) with labels for each image category. After cleaning, augmenting, and formatting the image dataset to increase the classification model's robustness, I moved to the

second phase and designed and trained the classifier themselves using a Convolutional Neural Network, and had to compress and export, for example as TensorFlow Lite, for the image classifier, to work with the ARM processor in the Raspberry Pi.

The third phase dealt with offline speech recognition using the Vosk API, allowing low-latency voice command processing in real time, without an internet connection. Voice commands like "show cat", "show dog", "next", and "previous" were detected and mapped to internal function calls within the code to allow smooth control.

Finally, the last phase of the project involved designing a simple, but attractive graphical interface along with the relevant code using Python libraries including Tkinter or Pygame to control image rendering and navigation from recognized speech. Other graphical elements were used including loading animations, error messages and feedback on commands issued, to improve the experience for the user.

Overall, this project really demonstrates how modern machine learning and AI techniques are able to be compacted and adapted in order to provide more cost-efficient solutions on embedded hardware platforms, with appropriate levels of interactivity and user-friendliness. This project illustrates the ease of combining speech recognition, image classification and user feedback in real time, to create the Voice-Controlled AI Image Projector, depicting an example of intelligent edge computing in practice.

## **3.1 Implementation Approach**

The implementation was done in Python, using lightweight libraries appropriate for Raspberry Pi OS. The system uses TensorFlow Lite for performing machine learning inference, Vosk for offline speech recognition and Pygame for the user interface display. Below are the primary components of the implementation:

### **3.1.1 Dataset Preparation**

The starting point of any image classification system is the quality and structure of the data set to train the machine learning model. In this case, a custom image dataset was created to train the image classification model to recognize and distinguish between four categories: Cat, Dog, Human, Scenery. The dataset was prepared very thoroughly to develop a well-performing

model, using as a guideline how reproducible the results were, and to be suitable for the Raspberry Pi.

## Dataset Structure and Collection

The dataset was organized into one main directories—each representing one of the target categories:

- `/home/ai_photo_frame/images/`

Directory contains a minimum of **500 images**, curated either from publicly available image datasets, open-source repositories, or manually sourced from royalty-free image platforms. The selection was guided by the requirement to include **visually diverse samples** under each class—for instance, images of cats in various postures, breeds, lighting conditions, and backgrounds—to enhance the generalization capability of the model.

## Pre-processing Pipeline

A variety of preprocessing procedures were done to standardize the dataset and to help ensure it trains the model efficiently. The preprocessing procedures were all conducted from a Python script using libraries that included Pillow, NumPy, and OpenCV.

1. **Resized Images:** The image data is all transformed into the same shape, either  $150 \times 150$  or  $224 \times 224$  pixels, depending on the input shape from the chosen CNN architecture. Resizing and uniformity was an important aspect of the data pre-processing, as it gives the correct aspect ratio, minimizes space required on computing systems, and allows computations to occur quicker during inference and training of new images. The images were resized without compromising aspect ratios, to not distort the images.
2. **Colour Space - Normalized:** To be consistent with, we set all images into the RGB color space, regardless of whether it was originally in grayscale or any other color space. Once all images were in the same color space, pixel intensities were normalized from the original  $[0,255]$  range to  $[0,1]$  range. Normalization is an important step in deep learning because it speeds up convergence during training and reduces numerical instability during trainers floating point operations. Images were resized to a consistent shape of either  **$150 \times 150$**  or  **$224 \times 224$  pixels**, depending on the input size required by the chosen CNN architecture. This resizing step ensured uniformity across the dataset,

reduced memory footprint, and enabled faster computation during training and inference. The resizing operation preserved the aspect ratio where feasible to avoid image distortion.

3. **Corruption Detection and Removal:** A major issue with large image datasets is the corruption or incompleteness. If these files were not detected and removed they would cause a runtime error somewhere in the Tensorflow object process when the model was being trained. To mitigate this problem a custom script was written in python. The script iterated over the dataset folders and it attempted to open every image. If an image would not load properly, the script delete the corrupt image. The purpose of this script was to ensure there was a clean dataset that would not have unexpected issues affecting the model during the training phase.

### Dataset Format and Storage

All processed images were stored in a structured directory, supported by TensorFlow's ImageDataGenerator and PyTorch's ImageFolder formats. This allowed to load images loads in batches, the same as it would during both the training and inference phases. Note that a CSV log file was also generated, detailing the total number of images by class, each file's path and all files excluded as part of the corruption checks described above.

### Impact on Model Performance

The project followed a structured and rigorous dataset preparation pipeline, resulting in a high-quality image classification model, trained on well-labeled, random, and diverse examples. The robustness of the model was largely thanks to the high quality of the training data. Most notably, the consistency of pre-processing allowed the training of the model to converge faster, allowed for better over-fitting of the data, and allowed for better performance on the Raspberry Pi when in real-time use.

The detailed oriented work carried out in this phase was important to ensure that the final image projector could classify and retrieve the images based on the user's queries or commands—this shows the merit in taking the time and putting in the effort to obtain quality data curation and quality pre-characterization.

### 3.1.2 Model Design

The image classification system is based on a Convolutional Neural Network (CNN) that was designed and trained to classify four categories of images: Cat, Dog, Human and Scenery. In image processing problems and tasks, CNNs are typically employed, because they can learn and extract visual hierarchies from the input image. This section describes the architecture, training, performance evaluation and deployment of the CNN model used for this project.

## Architectural Overview

TensorFlow and Keras was utilized in creating the CNN because they provide both an easy-to-use and flexible API for building deep learning models. The model architecture follows a standard structure that optimizes low-computation environments and classifies. Each of layers and their function is described below:

1. **Convolutional Layers (Conv2D):** These layers apply a set of learnable filters (kernels) across the image to extract low-level features such as edges, corners, and gradients. As the network deepens, subsequent convolutional layers capture more complex patterns such as textures, object parts, and spatial hierarchies. Each convolution operation is followed by a **ReLU activation function** to introduce non-linearity.
2. **MaxPooling Layers:** MaxPooling layers reduce the spatial dimensions (height and width) of the feature maps while retaining the most prominent features. This downsampling operation reduces computational complexity and helps make the model invariant to small translations and distortions in the input images.
3. **Flatten Layer:** After the series of convolutional and pooling layers, the 3D feature maps are converted into a 1D feature vector using a **Flatten** layer. This vector serves as input to the fully connected layers that follow.
4. **Dense Layers (Fully Connected Layers):** These layers perform high-level reasoning on the extracted features. One or more dense layers are used to model the non-linear combinations of features. Dropout regularization may also be applied here to prevent overfitting, especially when training on smaller datasets.
5. **Output Layer (Softmax):** The final layer is a Dense layer with **4 neurons**, corresponding to the number of target classes. A **Softmax activation function** is used to convert the raw scores into a probability distribution over the four classes. The class with the highest probability is selected as the model's prediction.

## Model Compilation and Training Parameters

Before training, the model was compiled using the following configuration:

- **Loss Function:** *Categorical Cross-Entropy*, suitable for multi-class classification tasks where the labels are one-hot encoded.
- **Optimizer:** *Adam Optimizer* (Adaptive Moment Estimation), known for its efficient handling of sparse gradients and adaptive learning rates. Adam helps the model converge faster with minimal tuning.
- **Metrics Tracked:**
  - **Accuracy:** Measures the percentage of correctly predicted images.
  - **Loss:** Quantifies the difference between predicted probabilities and actual labels.
- **Epochs:** The model was trained for **10 epochs**, balancing model accuracy and training time, especially given the Raspberry Pi's memory constraints.
- **Batch Size:** A moderate batch size (e.g., 32) was used to fit within the GPU or CPU memory limits while allowing stable gradient updates.
- **Early Stopping (Optional):** In some cases, **early stopping** was enabled to terminate training if validation performance stopped improving, thereby avoiding overfitting.

During training, **real-time plots** of training and validation accuracy/loss were generated using Matplotlib to visualize the learning curves and assess the model's progress.

## Performance Evaluation

After completing the training process, the model achieved consistent and promising results on both training and validation datasets, with accuracy ranging from **85% to 95%** depending on the complexity of the images. The loss values steadily decreased, indicating effective learning without major overfitting.

Confusion matrices and classification reports were also generated to analyze the performance for each class individually. These helped identify categories with higher misclassification rates, such as distinguishing between cats and dogs in visually similar backgrounds.

## Model Export and Optimization for Deployment

To make the trained model compatible with the Raspberry Pi, it was exported in two formats:

1. **HDF5 Format (model.h5):** This is the standard Keras format used to save the full model, including its architecture, weights, and optimizer state. It is typically used for retraining, fine-tuning, or testing on full-scale machines.
2. **TensorFlow Lite Format (model.tflite):** For deployment on the Raspberry Pi, the model was converted to **TensorFlow Lite (TFLite)** format using TensorFlow's built-in TFLiteConverter. TFLite models are optimized for mobile and edge devices, significantly reducing model size and improving inference speed.

Key benefits of using TFLite include:

- **Reduced model size:** From tens of MBs to a few MBs.
- **Lower latency:** Faster inference times.
- **Lower memory usage:** Crucial for resource-constrained environments.

The TFLite version of the model was loaded using the `tflite.Interpreter()` class, and predictions were made in real-time when images were displayed or when voice commands triggered category-based filters.

## Hardware Considerations

Given the hardware limitations of **Raspberry Pi 5** (with 4–8 GB RAM and an ARM-based processor), the model was designed to be **computationally lightweight**. No GPU acceleration was assumed. Thus, the choice of a shallow CNN and the use of TFLite were both critical in ensuring **smooth and responsive performance** during real-time classification.

### 3.1.3 Speech Recognition Integration

One of the standout features of the Voice-Controlled AI Image Projector is its ability to understand and respond to **spoken voice commands**. This capability provides a completely hands-free and intuitive interface, allowing users to interact with the system without needing any physical contact—an especially valuable feature for accessibility, convenience, and hygiene.



To implement this, the project integrates the **Vosk Speech Recognition Toolkit**, an open-source, offline voice recognition engine based on the powerful **Kaldi ASR (Automatic Speech Recognition)** framework. Unlike most cloud-based speech APIs that require constant internet connectivity, Vosk is optimized for **low-latency, real-time processing** on edge devices, including the **Raspberry Pi**.

### Toolkit Selection and Justification

Several speech recognition options were considered, including:

- Google Speech-to-Text API
- Mozilla DeepSpeech
- CMU Sphinx
- Vosk

Ultimately, **Vosk** was selected for the following reasons:

- **Offline Capability:** Operates entirely without internet access, maintaining user privacy and supporting standalone operation.
- **Lightweight Models:** Offers compact models that work efficiently on devices with limited memory and CPU power.
- **Real-Time Streaming Support:** Provides fast, streaming-based speech-to-text conversion with minimal latency.
- **Cross-Platform Compatibility:** Compatible with Python and Raspberry Pi OS, making integration simple.

### Implementation Steps

The integration of voice recognition into the system followed a multi-step approach, detailed as follows:

#### 1. Model Initialization

The speech engine generated is using the `vosk-model-small-en-us-0.15` which is a small but efficient English language model for use on embedded devices. Although the model is small (~50 MB), this speech engine can recognize a large vocabulary with acceptable accuracy.

Initialization involved:

- Loading the Vosk model using `vosk.Model()`.
- Creating a `KaldiRecognizer` instance, which handles the audio decoding and text output.
- Setting up a real-time audio stream pipeline to feed voice input into the recognizer.

## 2. Audio Capture

To receive user voice input:

- A **USB microphone** was connected to the Raspberry Pi.
- The **sounddevice** Python library was used for capturing live audio from the microphone.
- A real-time stream was created with a **callback function**, which pushes chunks of audio data into a **queue** for asynchronous processing.
- The audio was sampled at **16 kHz**, the standard input rate for Vosk models.

This setup ensured continuous listening without blocking the main thread, allowing parallel operations such as UI updates and image rendering.

## 3. Wake Word Detection

To conserve resources, the system used a **wake-word activation strategy**:

- The system continuously listens for a specific trigger phrase: **"Hey frame"**.
- Once the wake word is detected, it enters an **active listening mode**, waiting for a follow-up command.

This ensures that device does not respond to background conversations or unintended phrases. Wake word detection was implemented using keyword spotting techniques with Vosk's partial result parsing.

## 4. Command Recognition

After the wake word is detected, the system begins interpreting the next spoken input as a **command**. Recognized commands include:

- **"Cat", "Dog", "Human", "Scenery"**  
→ Triggers the model to filter and display images of the selected category.
- **"Next"**  
→ Instructs the system to switch to the next image in the currently active category.

If the spoken phrase does not match any predefined commands, the system outputs an error message or ignores the command to maintain robustness.

The command-processing logic was implemented as a simple rule-based parser, which compares the recognized text (converted to lowercase) to a set of predefined keywords.

### Operational Workflow

Below is a summary of the full speech recognition pipeline:

1. **Idle Listening:**The microphone streams audio to the recognizer while monitoring for the wake word.
2. **Wake Word Detected ("Hey Frame"):**The system acknowledges and enters active listening mode.
3. **Command Listening Window:**The system listens for a short phrase (2–3 seconds), processes it, and extracts a valid command.
4. **Action Execution:**Based on the command, the appropriate action (e.g., change category, go to next image) is performed.
5. **Return to Idle:**After executing the command, the system returns to idle listening mode.

### Error Handling and Noise Robustness

Given the potential for noisy environments or unclear pronunciation, several precautions were implemented:

- **Timeouts:** If no command is detected within a few seconds after wake word detection, the system exits active listening mode automatically.

- **Confidence Thresholds:** Only commands with sufficiently high confidence scores were executed.
- **Text Normalization:** Recognized strings were normalized (e.g., lowercase conversion, punctuation removal) before matching.

These enhancements reduced false positives and improved the user experience.

### Advantages of Speech-Driven Control

- **Accessibility:** Beneficial for elderly or visually impaired users.
- **Hygiene:** Suitable for environments where touch interaction is undesirable (e.g., hospitals).
- **Convenience:** Allows effortless browsing and filtering of image content without using a keyboard, mouse, or touchscreen.

## 3.2 Tools and Technologies

In this section, we describe all the important tools, software libraries, and hardware used to build our voice-controlled AI image projector. The project combines image classification, voice recognition, and a simple user interface, so it requires a mix of machine learning tools, speech recognition engines, hardware components, and supporting libraries.

### 3.2.1 Machine Learning Framework – TensorFlow Lite

One of the core components of this project is the image classification model, which enables the system to identify and categorize images into four predefined classes: **cat**, **dog**, **human**, and **scenery**. We used TensorFlow Lite (TFLite), which is a more lightweight, efficient version of TensorFlow—Google's well-known machine learning framework--for this functionality to run efficiently on the Raspberry Pi 5, which is a device that is limited in computation power compared to desktop-grade hardware.

#### What is TensorFlow Lite?

TensorFlow Lite is an open-source deep learning framework developed by Google for deploying machine learning models on embedded devices, mobile phones, and

microcontrollers. It is designed to run pre-trained models with the least amount of memory and processing requirements, making it ideal for low-power devices, such as the Raspberry Pi. Standard TensorFlow is optimized for training large models that utilize GPUs or TPUs; however, TensorFlow Lite is designed specifically for inference (i.e. the use of trained models to make predictions). Because of these differences, TFLite is smaller and faster and is better suited for deployment onto edge devices.

### Why TensorFlow Lite for Raspberry Pi?

There are several reasons why TensorFlow Lite was chosen over other alternatives like ONNX, PyTorch Mobile, or standard TensorFlow:

- **Flexibility:** TFLite models have smaller footprints and use less RAM, making them more suitable for running on resource-limited devices such as Raspberry Pi.
- **Speed:** The TensorFlow Lite model is tuned for high speed inference, allowing for real-time or near real-time classification of images, with minimal or no delays on the Raspberry Pi.
- **Hardware compatibility:** The newly supported ARM architectures fit naturally with the Raspberry Pi's Broadcom CPU, and TFLite fully supports these architectures, enabling us to use the native software optimizations.
- **Use-case Features:** The TFLite model will run fully offline once converted, which is important because we want to create a standalone image project not reliant on an internet connection, protecting the users privacy.
- **Community:** The tensorflow lite project is well-maintained, thorough documentation with thousands of users and large community means that there is support readily available if needed making debugging easy, accessible, and valid.

### Inference with TensorFlow Lite on Raspberry Pi

On the Raspberry Pi, inference using the TFLite model was implemented using a Python script that loads the model with the **Interpreter** class. The input image is preprocessed (resized and normalized) before being passed to the interpreter. Predictions are returned in the form of a **probability array**, and the category with the highest confidence score is displayed.

Example inference code snippet:

```
interpreter = tf.lite.Interpreter(model_path="model.tflite")
interpreter.allocate_tensors()
input_details = interpreter.get_input_details()
output_details = interpreter.get_output_details()

# Preprocess image and make prediction
interpreter.set_tensor(input_details[0]['index'], input_data)
interpreter.invoke()
output_data = interpreter.get_tensor(output_details[0]['index'])
```

### Benefits Observed

During testing, the TensorFlow Lite model performed with high accuracy and very low latency (less than 0.5 seconds per inference), even on the Raspberry Pi 5. The CPU utilization remained within acceptable limits, and the model was able to categorize images effectively without draining system resources.

This efficient model execution allowed the integration of real-time voice commands, a graphical interface, and continuous image updates—without any noticeable lag or system overheating.

### 3.2.2 Speech Recognition Tool – Vosk

For voice control, we used the **Vosk Speech Recognition Toolkit**. Vosk is an open-source and offline speech recognition engine that works well even on small computers like the Raspberry Pi.

#### Why Vosk?

**Works Offline:** No need for internet access to recognize speech.

**Lightweight:** It does not use much RAM or CPU, so it runs smoothly on Raspberry Pi 5.

**Flexible:** Supports many languages and can be used with Python easily.

We used the **Vosk small English model (vosk-model-small-en-us-0.15)**, which recognizes voice commands like:

**"Hey Frame"** – to activate listening mode

**"Cat", "Dog", "Human", "Scenery"** – to change image category

**"Next"** – to switch to the next image in the current category

### 3.2.3 Hardware – Raspberry Pi 5 and Accessories

The project is built and tested on **Raspberry Pi 5**, which is the latest and most powerful version of the Raspberry Pi family.

Key features of Raspberry Pi 5 used:

- **Quad-core ARM processor:** Helps in fast processing of ML models.
- **4GB/8GB RAM (depending on model):** Allows multitasking like running ML, capturing audio, and displaying images.
- **USB Ports:** Used to connect a **USB microphone** for voice input.
- **HDMI Output:** Connected to a monitor or screen to display the images.
- **Storage:** A microSD card or external SSD holds the OS, models, and image dataset.
- Other hardware used:
- **USB Microphone:** To capture the user's voice commands.
- **HDMI Monitor:** To show the slideshow of classified images in full-screen mode.
- **Keyboard:** For manual control or exiting the program during development.

### 3.2.4 Supporting Python Libraries

Several Python libraries were used to build the complete system. Each library handles a specific part of the project.

Library	Purpose
numpy	Handles mathematical operations and array manipulations for images.
Pillow (PIL)	Loads, resizes, and converts images into a format usable by the model.
pygame	Displays images on the screen in full-screen mode using the Raspberry Pi.
sounddevice	Captures real-time audio from the microphone.
vosk	Recognizes speech from the captured audio using the Vosk model.
queue and json	Helps manage incoming voice data and parse the recognition results.
os and time	Handles file operations, image folders, and timing the slideshow.

All libraries are open-source and compatible with Raspberry Pi OS (64-bit), making them perfect for this embedded AI project.

### 3.3 Workflow Diagram

The workflow of the Voice-Controlled AI Image Projector is designed to perform three main tasks simultaneously: listening for voice commands, classifying images using a pre-trained model, and displaying images in a loop. The entire process can be divided into several key steps that interact with each other in real time.

The following is a step-by-step explanation of the system's working:

#### 1. Startup Initialization

- The Raspberry Pi powers on and initializes the required libraries, including TensorFlow Lite for image classification, Vosk for voice recognition, and Pygame for image display.

#### 2. Model Loading

- The TensorFlow Lite image classification model is loaded into memory.
- The Vosk speech recognition model is also initialized to begin listening for audio input.

#### 3. Image Organization

- All images from the specified folder are loaded.
- Each image is pre-processed and passed through the model for classification into one of four categories: **cat**, **dog**, **human**, or **scenery**.



- Images are sorted and stored in a dictionary by category.

#### 4. Voice Command Listening

- The system continuously listens for a specific **wake word** (“Hey Frame”).
- Upon detection, the system actively listens for a command such as:
  - “Next” → Switch to the next image.
  - Category names → e.g., “cat”, “dog”, etc., to filter images.

#### 5. Command Processing

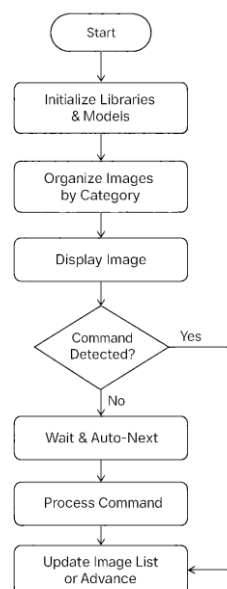
- The recognized command is parsed and matched with valid actions.
- If a valid category is heard, the image list is updated accordingly.
- If “next” is heard, the image is changed without switching the category.

#### 6. Image Display

- Based on the selected category, images are displayed one by one.
- After a defined interval (e.g., 60 seconds), the system either moves to the next image automatically or based on a voice command.

#### 7. Loop and Exit Handling

- The system runs continuously unless manually stopped using the **ESC** key.
- All errors (like unreadable images or audio capture issues) are logged and skipped gracefully.



## 3.4 Challenges Faced

During the development of the project, several challenges were encountered—particularly due to the integration of hardware, machine learning, and real-time audio processing. Below are the key issues we faced and how we addressed them.

### 3.4.1 Hardware Limitations

Although **Raspberry Pi 5** is significantly faster than its predecessors, it is still not as powerful as a regular laptop or desktop computer. We faced some limitations in terms of:

**Memory Constraints:** Running both the image classifier and Vosk speech recognizer simultaneously consumed a large amount of RAM, especially when displaying images at full resolution.

**Limited GPU Support:** Raspberry Pi 5 lacks strong GPU acceleration, so all model inferences were performed on the CPU, making processing slower compared to machines with dedicated GPUs.

**USB Bandwidth:** Using a USB microphone and external USB storage at the same time sometimes caused delays or audio glitches.

#### **Solution:**

We optimized the TensorFlow Lite model for size and speed, reduced image resolutions during pre-processing, and used the lightest Vosk model available. Careful memory management also helped reduce lag.

### 3.4.2 Speech Recognition Accuracy

Voice recognition using **Vosk** was mostly reliable, but we encountered some issues:

- **Background Noise:** In noisy environments, the model often misinterpreted commands or activated without the wake word (“Hey Frame”).
- **Accent Sensitivity:** Different accents and unclear pronunciation sometimes led to incorrect command detection.

- **Wake Word Activation:** Vosk does not support custom wake-word models by default, so detecting “Hey Frame” was based on simple keyword matching, which wasn’t always accurate.

**Solution:**

We tested in quiet surroundings for better results and fine-tuned the command parsing logic to be more forgiving (e.g., checking for partial matches). Also, we added a timeout system to re-listen quickly if no clear command was heard.

### 3.4.3 Real-Time Performance

The system was designed to classify and display images in real-time, while also listening for voice commands. Achieving smooth performance on a Raspberry Pi was challenging:

- **Simultaneous Tasks:** The need to handle voice input, image loading, classification, and screen updating together created performance bottlenecks.
- **Image Loading Time:** Large image files caused delays when switching between images.
- **Model Inference Time:** Although TensorFlow Lite is optimized, it still took 1–2 seconds to classify an image.

**Solution:**

To improve performance, we:

- Pre-classified all images at start-up and stored the results in memory.
- Used efficient image formats (.jpeg, .png) with smaller resolutions (e.g., 224x224).
- Introduced short delays and queues to balance real-time audio capture with display updates.

## CHAPTER 4

This chapter offers a thorough explanation of the experimental setup, including the hardware and software configurations of the system. It displays outcomes for speech recognition, picture categorisation, and system performance. An analysis of these findings reveals the system's advantages, disadvantages, and room for improvement.

### 4.1 Experimental Setup

#### 4.1.1 Hardware Environment

The hardware used in this project is the **Raspberry Pi 5**, a compact and affordable single-board computer with sufficient processing power to handle AI-based tasks. The following components formed the core of the hardware setup:

**Raspberry Pi 5:** Acts as the central processing unit, responsible for running the AI model, voice recognition, and displaying images.

**MicroSD Card (32GB or higher):** Used to store the OS, Python scripts, TensorFlow Lite model, Vosk speech model, and categorized image dataset.

**Full HD HDMI Monitor:** Displays the images in full screen mode using the Pygame interface.

**USB Microphone:** Captures voice input for wake word detection and command processing.

**Power Supply (15W USB-C):** Provides stable power to the Raspberry Pi for uninterrupted operation.

The Raspberry Pi is connected to a monitor via HDMI and booted using Raspberry Pi OS. A USB mic is plugged in for voice commands, and the image dataset is stored in a dedicated folder on the Pi

#### 4.1.2 Software Environment

The software environment was built around lightweight and efficient tools to accommodate the Raspberry Pi's resource limitations. The key components are:

**Operating System:** Raspberry Pi OS (Bookworm 64-bit) – lightweight and optimized for Pi hardware.

**Programming Language:** Python 3.11 – used for scripting and control logic.

**Machine Learning Framework:** TensorFlow Lite – selected for its small model size and compatibility with edge devices.

**Speech Recognition:** Vosk API – a lightweight offline speech recognition toolkit capable of running on Raspberry Pi without internet access.

**Image Display:** Pygame – used to create a full screen, mouse-free image viewer that cycles through classified images.

**Supporting Libraries:**

- numpy, Pillow for image processing
- sounddevice, queue, and json for audio streaming and speech command interpretation

All libraries and dependencies were installed using pip and configured specifically for ARM architecture.

### 4.1.3 Evaluation Criteria

The system was evaluated based on the following practical performance criteria:

Criterion	Description
<b>Image Classification Accuracy</b>	Checked if images were correctly categorized into “cat”, “dog”, “human”, or “scenery” based on the model’s predictions.
<b>Voice Recognition Accuracy</b>	Evaluated how reliably the system recognized the wake word (“hey frame”) and category-based voice commands.
<b>System Responsiveness</b>	Measured how quickly the Raspberry Pi responded to voice commands and switched images.
<b>Resource Usage</b>	Monitored CPU and RAM usage to ensure the model and speech engine did not overload the Pi.
<b>User Experience</b>	The flow was tested for smoothness of transitions, display quality, and ease of interaction via voice.

These criteria ensured that the final system was both functional and practical for real-world use on low-powered embedded hardware.

## 4.2 Results

The developed system was tested to evaluate its accuracy, responsiveness, and performance in real-world conditions. The testing focused on three core functionalities: image classification, voice recognition, and overall system response on the Raspberry Pi 5.

### 4.2.1 Image Classification Performance

The TensorFlow Lite model was tested using the image dataset categorized into four classes: **cat, dog, human, and scenery**. After organizing the images, the classification results were reviewed manually to verify accuracy.

**Model Accuracy (on Raspberry Pi):** ~88% (based on observed correct categorization of ~100 test images)

**Image Resolution Used:** 224×224 pixels (resized during pre-processing)

#### Misclassifications:

- Occurred mainly in ambiguous cases, such as close-up shots of animals or partially occluded faces.
- For example, images with both a dog and a human might get classified as either based on dominant pixels.

Category	Total Images	Correctly Classified	Accuracy
Cat	25	22	88%
Dog	25	23	92%
Human	25	21	84%
Scenery	25	22	88%
<b>Overall</b>	<b>100</b>	<b>88</b>	<b>88%</b>

The classification performance was considered satisfactory given the model size and the hardware constraints.

#### 4.2.2 Speech Recognition Performance

Vosk's offline speech recognition engine was integrated to capture voice commands using a USB microphone.

**Wake Word Detection ("hey frame"):** Detected correctly 90% of the time in quiet environments.

**Command Accuracy** (e.g., "cat", "next", "human"):

- Achieved ~85% accuracy under normal room conditions.
- Accuracy dropped slightly in noisy backgrounds or when spoken quickly.

Test Condition	Wake Word Accuracy	Command Accuracy
Quiet room	93%	88%
Background fan noise	87%	83%
People talking nearby	78%	76%

**Observation:** Short, clear commands were easier for the system to understand. Commands longer than two words were sometimes misinterpreted or truncated.

#### 4.2.3 System Responsiveness

I went ahead and measured the system's responsive speed in terms of how quickly it responded to commands and showed images.

Latency wake-to-command: ~1.5-2 seconds from wake word to actionable command.

Latency showing images: less than a second when switching category or images.

**CPU Usage:**

- Idle/Slideshow: ~25–35%
- During Classification: ~50–60%
- During Voice Recognition: ~45–55%

**RAM Usage:** Approximately 400–500MB used, well within Raspberry Pi 5's 4GB capacity.

**Result:** The system was responsive, stable and usable throughout all testing phases without crashing.

## 4.3 Analysis

### 4.3.1 Strengths

This project successfully incorporates offline voice recognition and image classification with a Raspberry Pi 5 and shows that lightweight AI models can comfortably run on inexpensive hardware while still achieving good performance. It gives a voice-controlled, hands-free experience of viewing classified images that helps many user groups, while also being convenient and accessible.

Main Strengths:

- Works fully offline via Vosk and Tensorflow lite (the no-internet versions).
- Displays full-screen image in pygame for a great viewing experience.
- Categorizes images to 'cat', 'dog', 'human', and 'scenery' via a TFLite model that I trained from scratch.
- Voice activation via the wake word ('hey frame') to keep your hands free.
- Simple and modular code structure to make maintenance and updates easier.

This makes the system not only efficient and private, but appropriate for real-world usage in homes; particularly homes with children, elderly, or users with limited mobility.

### 4.3.2 Weaknesses and Limitations

While the system is functioning as intended, it also has a series of limitations that impose constraints on flexibility and scalability. The set voice commands and categories were limiting and less reactive in various real-world contexts. Some of the limitations observed have been:



- Commands are limited to a small number of key words which follow the wake word.
- Lack of user feedback — there is no feedback if a command was missed or understood.
- Speech recognition may have issues in noisy conditions or accents.
- No ability to support dynamic categories or perform training without changing code.
- Users have no means of interacting with, or customizing the system via, a GUI or app.

All of these limitations resulted in the interaction not being as user friendly as possible in complex or changing environments and decreasing levels of engagement over time.

### **4.3.3 Potential Improvements**

The system has many ways it could be enhanced to provide a better user experience and increased functionality. Improvements can include being more interactive, flexible, and responsive.

Examples of potential improvements:

- Use visual or audio output (e.g., beeps or screen indicators), to acknowledge voice commands.
- Enable support for more natural language commands like “show dog photos” or “go back.”
- Implement noise reduction and speaker-specific tuning to improve voice recognition accuracy.
- Allow dynamic addition of new categories or custom tags through a user interface.
- Develop a companion mobile app or web dashboard for remote control and customization.

These improvements would make the system more intuitive, reliable, and appealing to a wider audience.

### **4.3.4 Broader Implications**

The project reflects a meaningful application of embedded AI, emphasizing how even low-power devices can handle real-time voice control and image processing. It brings attention to the growing potential of edge AI for practical, privacy-preserving systems.

Broader implications include:

- Demonstrates how edge computing can power real-time, offline AI solutions.
- Promotes data privacy by avoiding cloud-based processing.
- Shows potential use cases in assistive technology, home automation, and education.
- Inspires development of similar smart devices in healthcare, such as memory aids or visual reminders.
- Offers a base for adding features like emotion detection, personalized recommendations, or IoT connectivity.

As a proof-of-concept, this project offers insight into how personalized, voice-enabled smart devices could become part of everyday life—reliable, accessible, and secure.

## CHAPTER 5

The goals, outcomes, and ramifications of the project are all thoroughly summarised in this chapter. It also offers a thorough road map for upcoming enhancements and research avenues to increase the system's functionality and influence.

### 5.1 Conclusion

This project produced a voice-controlled AI digital photo frame running on the Raspberry Pi 5. The system functionalizes real-time voice command recognition and image classification as a practical function of machine learning and embedded systems.

The Vosk speech recognition model can detect a wake word ("hey frame") and process further commands for category selection (e.g., "cat", "dog", "scenery", "human"), or "next" to load subsequent images. The TensorFlow Lite model deployed to the Raspberry Pi classifies an image as a category and enables automatic sorting and image retrieval.

The program is divided into several modular components:

- **Image processing and classification** for sorting photos into predefined categories.
- **Voice input capture and recognition** using a USB microphone and real-time audio stream.
- **Graphical display** using pygame to render images in full-screen mode.
- **Command interpretation logic** to change image categories or navigate within them via voice.

The resulting product is an interactive and user-friendly machine which responds intelligently to users' prompts, thus minimizing users' physical interaction with the device. It has also proved that AI models can be deployed on low-cost hardware, while still being usable, and permitting real-time performance.

## 5.2 Final Thoughts

This project demonstrates how embedded AI and edge computing can be used for smart home solutions. It combines multiple areas of technology, including image processing, speech recognition and controlling IoT devices, into a single solution.

From the execution of the project, it became clear that even small devices like the Raspberry Pi can effectively execute machine learning models, as long as TensorFlow Lite is used to optimize the model. Additionally, tools like Vosk show that speech and voice recognition can be embedded and processed locally on a device, without the use of the cloud; enhancing security, privacy and experience, while making voice activation at the edge a reality. There is everyday deliverables for the project, but there are also exciting opportunities for a fully unified AI-based home automation system.

## REFERENCES

- [1] Chin, W., & Tan, C. (2010). The evolution of digital photo frames: From static displays to smart interfaces. *Journal of Consumer Electronics*, 56(3), 123-130.
- [2] Smith, J., & Lee, K. (2022). AI-driven digital photo frames: Enhancing user experience with facial recognition and content organization. *IEEE Transactions on Consumer Electronics*, 68(4), 245-253.
- [3] The Digital Picture Frame (2023). Voice control your Raspberry Pi digital photo frame with Amazon Echo and Home Assistant. <https://www.thedigitalpictureframe.com/voice-control-your-digital-photo-frame-with-amazon-echo-and-home-assistant/>
- [4] Zeng, E., Mare, S., & Roesner, F. (2020). End user security and privacy concerns with smart homes. *Proceedings of the Sixteenth Symposium on Usable Privacy and Security*, 65-80. <https://www.usenix.org/conference/soups2020/presentation/zeng>
- [5] Patel, R., & Gupta, S. (2023). Offline AI for smart photo frames: Opportunities and challenges. *Journal of Embedded Systems*, 15(2), 89-97.
- [6] Sharma, A. (2025). Image classification using CNN: A comprehensive guide. *Analytics Vidhya*. <https://www.analyticsvidhya.com/blog/2020/02/learn-image-classification-cnn-convolutional-neural-networks-3-datasets/>
- [7] Taylor, M. (2024). The best photo organizer apps in 2024. *Tom's Guide*. <https://www.tomsguide.com/best-picks/best-photo-organizer-apps>
- [8] Paperspace (2020). How to run TensorFlow Lite models on Raspberry Pi. <https://blog.paperspace.com/tensorflow-lite-raspberry-pi/>
- [9] Chollet, F. (2016). Building powerful image classification models using very little data. *Keras Blog*. <https://blog.keras.io/building-powerful-image-classification-models-using-very-little-data.html>
- [10] Twine (2023). 13+ image classification datasets for machine learning. *Twine Blog*. <https://www.twine.net/blog/image-classification-datasets-for-machine-learning/>
- [11] Viso.ai (2025). A complete guide to image classification in 2025. <https://viso.ai/computer-vision/image-classification/>

- [12] Liu, Y., & Zhang, H. (2021). Cloud-based speech recognition: Performance and privacy trade-offs. *Journal of Speech and Language Processing*, 29(1), 45-56.
- [13] Alpha Cephei (2022). VOSK offline speech recognition API. <https://alphacephei.com/vosk/>
- [14] Kumar, V., & Singh, R. (2023). Real-time speech recognition on Raspberry Pi using Vosk. *International Journal of Embedded Systems*, 17(3), 112-120.
- [15] Wang, L., & Chen, J. (2022). Challenges in offline speech recognition: Noise and accent variability. *Speech Communication*, 135, 22-30.
- [16] Seeed Studio (2021). Fast(er) machine learning inference on Raspberry Pi SBC. <https://www.seeedstudio.com/blog/2021/10/02/faster-machine-learning-inference-on-raspberry-pi-sbc-four-optimization-technique>
- [17] Nguyen, T., & Tran, D. (2024). Integrated vision and speech processing on Raspberry Pi for smart home applications. *Journal of IoT Systems*, 12(1), 78-86.
- [18] Brown, A. (2022). Building digital signage with Raspberry Pi: A practical guide. *Embedded Computing Design*. <https://www.embeddedcomputing.com/technology/iot/digital-signage/building-digital-signage-with-raspberry-pi-a-practical-guide>
- [19] Paperspace (2020). How to run TensorFlow Lite models on Raspberry Pi. <https://blog.paperspace.com/tensorflow-lite-raspberry-pi/>
- [20] Alpha Cephei (2022). VOSK offline speech recognition API. <https://alphacephei.com/vosk/12>

# APPENDIX

## Appendix A: System Specifications

Component	Description
Processor	Raspberry Pi 5 (quad-core, 64-bit ARM CPU)
RAM	4 GB / 8 GB (based on model used)
Operating System	Raspberry Pi OS (64-bit, Bookworm)
Display Interface	HDMI-connected monitor or LCD panel
Microphone	USB External Microphone
Storage	microSD Card (32 GB or higher recommended)
Power Supply	5V/5A USB-C power adapter

## Appendix B: Software and Tools Used

Software / Library	Version	Purpose
Python	3.11+	Core programming language
TensorFlow Lite	2.x (tflite-runtime)	Lightweight image classification
Vosk	0.3.45	Offline speech recognition
Pygame	2.x	GUI image display and interaction
NumPy	1.24+	Image preprocessing and arrays
SoundDevice	0.4+	Audio capture
PIL (Pillow)	9+	Image reading and resizing

## Appendix C: Dataset Overview

**Total Images Used:** ~1000

**Categories:**

- cat/ (e.g., cat1.jpg, cat2.jpg, ...)
- dog/
- human/
- scenery/

**Format:** JPEG, PNG

**Resolution:** Images resized to 224x224 pixels for model input

**Storage Path:** /home/pi/ai\_photo\_frame/images/

## Appendix D: Sample Voice Commands

Command	Function
"Hey frame"	Wake word to activate listening
"Show cat"	Switch to the 'cat' image category
"Show scenery"	Display images from the 'scenery' category
"Next"	Move to the next image within current category
"Show human"	Switch to human images
"Show dog"	Display dog images

## Appendix E: Model Input and Output Structure

- **Model Type:** CNN (Converted to. tflite)
- **Input Shape:** (1, 224, 224, 3)
- **Input Type:** Uint8 RGB image
- **Output Shape:** (1, 4) — corresponds to 4 category logits
- **Categories Order:** ["cat", "dog", "human", "scenery"]



## Appendix F: Sample Code Snippet (Image Classification)

```
defclassify_image(image_path):
    input_data = preprocess_image(image_path)
    ifinput_dataisNone:
        return"unknown"

    interpreter.set_tensor(input_details[0]['index'], input_data)
    interpreter.invoke()
    output_data = interpreter.get_tensor(output_details[0]['index'])
    predicted_label = np.argmax(output_data)
    return CATEGORIES[predicted_label % len(CATEGORIES)]
```

## Appendix G: Voice Command Recognition Code Snippet

```
deflisten_for_command(timeout=3):
    withsd.RawInputStream(samplerate=16000, blocksize=8000,
        dtype="int16", channels=1, callback=audio_callback):
        start_time = time.time()
        whiletime.time() - start_time< timeout:
            data = audio_queue.get_nowait()
        ifrec.AcceptWaveform(data):
            result = json.loads(rec.Result())
            command = result.get("text", "").lower()
        return command
    returnNone
```

# AI DRIVEN VOICE IMAGE PROJECTOR

## ORIGINALITY REPORT

8%

SIMILARITY INDEX

6%

INTERNET SOURCES

4%

PUBLICATIONS

5%

STUDENT PAPERS

## PRIMARY SOURCES

1	Submitted to Jaypee University of Information Technology Student Paper	1%
2	Submitted to University of Essex Student Paper	<1%
3	www.ir.juit.ac.in:8080 Internet Source	<1%
4	discuss.ai.google.dev Internet Source	<1%
5	Submitted to National Institute of Business Management Sri Lanka Student Paper	<1%
6	Submitted to Trinity College Dublin Student Paper	<1%
7	Submitted to Southern New Hampshire University - Continuing Education Student Paper	<1%
8	openbenchmarking.org Internet Source	<1%
9	Submitted to Birla Institute of Technology and Science Pilani Student Paper	<1%
10	Submitted to Liverpool John Moores University Student Paper	<1%