

# Visvesvaraya Technological University

Jnana Sangama, Belagavi – 590018



Project report on

## **“IMAGE TO SPEECH CONVERTER USING MACHINE LEARNING”**

**(PID:2023P008)**

Submitted in partial fulfillment for the award of Degree of Bachelor of Engineering

Submitted by

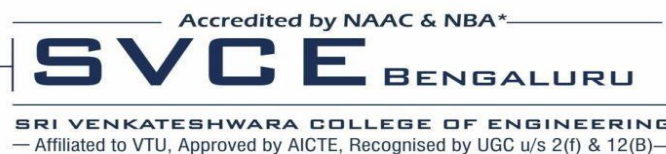
<b>ADITI SHARMA</b>	<b>1VE20CA001</b>
<b>GHANASHREE B N</b>	<b>1VE20CA008</b>
<b>MANOJ K</b>	<b>1VE20CS084</b>
<b>MOHITH GOWDA B M</b>	<b>1VE20CS090</b>

Under the Guidance of

**Prof. B PADMAVATHY**

Assistant professor

Dept. of CSE, SVCE, Bengaluru



**Department of Computer Science and Engineering**

**Computer Science Engineering in Artificial Intelligence**

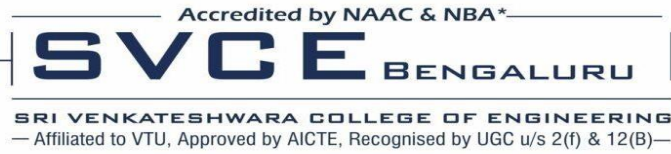
**Sri Venkateshwara College of Engineering**

**Bengaluru - 562157**

**ACY: 2023-24**

**Sri Venkateshwara College of Engineering, Bengaluru - 562157**

**Department of CSE&CSE-AI Engineering**



## **CERTIFICATE**

Certified that the Project work entitled **“IMAGE TO SPEECH CONVERTER USING MACHINE LEARNING”** was carried out by **ADITI SHARMA (1VE20CA001), GHANASHREE B N (1VE20CA008), MANOJ K (1VE20CS084), MOHITH GOWDA B M (1VE20CS090)**, who are bonafide students of VIII Semester, Computer Science and Engineering & CSE – Artificial Intelligence, Sri Venkateshwara College of Engineering, Bengaluru. This is in partial fulfillment for the award of Bachelor of Engineering in the Visvesvaraya Technological University, Belagavi, during the year **2023-24**. The project report has been approved as it satisfies the academic requirements in respect of Project work prescribed for the said Degree. To my best knowledge, the work reported in this thesis has not been submitted by us elsewhere for the award of the degree and is not the repetition of the work carried out by others.

.....	.....	.....	.....
<b>Signature of Guide</b>	<b>Signature of HOD</b>	<b>Signature of HOD</b>	<b>Signature of Principal</b>
<b>Prof. B Padmavathy</b>	<b>Dr. Hema M S</b>	<b>Dr. Prathima V R</b>	<b>Dr. Nageswara Guptha M</b>
Asst. Prof., Dept. of CSE	HOD, Dept. of CSE	HOD, Dept. of CSE-AI	Principal
SVCE, Bengaluru	SVCE, Bengaluru	SVCE, Bengaluru	SVCE, Bengaluru

Name of the examiners

1.

2.

Signature with Date:

## ACKNOWLEDGEMENT

The satisfaction that accompanies the successful completion of any task would be incomplete without mentioning the people whose ceaseless cooperation made it possible, whose constant guidance and encouragement crown all efforts with success.

Our sincere thanks to the highly esteemed institution **SRI VENKATESHWARA COLLEGE OF ENGINEERING**, for grooming us to be budding software engineers.

Our deepest thankfulness goes to our Principal, **Dr. NAGESWARA GUPTHA M** for his continuous encouragement and support.

Our deepest thankfulness goes to our Dean of Academics, **Dr. POORNIMA G R**, Department of Electronics and Communication Engineering. We do consider ourselves privileged to have had the opportunity to work under her guidance.

We would like to thank beloved **Dr. HEMA M S**, Head of the Department of Computer Science and Engineering and **Dr. PRATHIMA V R**, Head of the Department of Computer Science and Engineering – Artificial Intelligence respectively for their encouragement and support.

Our deepest thankfulness goes to our project coordinator and our internal guide, **Prof. B PADMAVATHY**, Assistant Professor, Department of Computer Science and Engineering. We do consider ourselves privileged to have had the opportunity to work under his guidance.

We thank our parents who were on our side whenever we needed them, and the almighty whose blessings are always there.

## **INSTITUTE VISION**

To be a premier institute for addressing the challenges in global perspective.

## **INSTITUTE MISSION**

**M-1:** Nurture students with professional and ethical outlook to identify needs, analyze, design and innovates sustainable solutions through lifelong learning in service of society as individual or a team.

**M-2:** Establish State-of-the Art Laboratories and Information Resource center for education and Research.

**M-3:** Collaborate with Industry, Government Organization and Society to align the curriculum and out reach activities.

## **DEPARTMENT VISION**

To get excellence in the field of artificial intelligence by providing technical education, research, innovation to meet the needs of industry and society.

## **DEPARTMENT MISSION**

- Mission 1: For acquiring the new skills and technology of artificial technology in software field and research field.
- Mission 2: To provide high quality of education and training in the field of artificial intelligence to get competence in areas of education by quality technical education and training.
- Mission 3: To get excellence in the technical education through education, research and service needed by the society.

## **PROGRAMME EDUCATIONAL OBJECTIVES (PEO's)**

**PSO1:** The ability to understand, analyse and demonstrate the knowledge of human cognition, Artificial Intelligence in terms of real world problems to meet the challenges of the future.

**PSO2:** The ability to develop computational knowledge and project development skills using innovative tools and techniques to solve problems in the areas related to Artificial Intelligence.



## **PROGRAM OUTCOMES**

Engineering Graduates will be able to:

- 1. Engineering knowledge:** Apply the knowledge of mathematics, science, engineering fundamentals, and an engineering specialization to the solution of complex engineering problems.
- 2. Problem analysis:** Identify, formulate, review research literature, and analyze complex engineering problems reaching substantiated conclusions using first principles of mathematics, natural sciences, and engineering sciences.
- 3. Design/development of solutions:** Design solutions for complex engineering problems and design system components or processes that meet the specified needs with appropriate consideration for the public health and safety, and the cultural, societal, and environmental considerations.
- 4. Conduct investigations of complex problems:** Use research-based knowledge and research methods including design of experiments, analysis and interpretation of data, and synthesis of the information to provide valid conclusions.
- 5. Modern tool usage:** Create, select, and apply appropriate techniques, resources, and modern engineering and IT tools including prediction and modeling to complex engineering activities with an understanding of the limitations.
- 6. The engineer and society:** Apply reasoning informed by the contextual knowledge to assess societal, health, safety, legal and cultural issues and the consequent responsibilities relevant to the professional engineering practice.

- 7. Environment and sustainability:** Understand the impact of the professional engineering solutions in societal and environmental contexts, and demonstrate the knowledge of, and need for sustainable development.
- 8. Ethics:** Apply ethical principles and commit to professional ethics and responsibilities and norms of the engineering practice.
- 9. Individual and team work:** Function effectively as an individual, and as a member or leader in diverse teams, and in multidisciplinary settings.
- 10. Communication:** Communicate effectively on complex engineering activities with the engineering community and with society at large, such as, being able to comprehend and write effective reports and design documentation, make effective presentations, and give and receive clear instructions.
- 11. Project management and finance:** Demonstrate knowledge and understanding of the engineering and management principles and apply these to one's own work, as a member and leader in a team, to manage projects and in multidisciplinary environments.
- 12. Life-long learning:** Recognize the need for and have the preparation and ability to engage in independent and life-long learning in the broadest context of technological change.

## ABSTRACT

we present a novel approach to image-to-speech conversion using machine learning techniques. Our system integrates a Generative Image-to-text Transformer (GIT) model, which unifies various vision-language tasks such as image/video captioning and question answering. Unlike previous methods that rely on complex structures and external modules, our approach simplifies the architecture into a single image encoder and text decoder, enhancing performance through scaled-up pre-training data and model size. Notably, our GIT model achieves state-of-the-art results on multiple benchmarks, including surpassing human performance on TextCaps. Additionally, we implement text-to-speech (TTS) functionality and provide camera access, enabling direct conversion of real images to speech. This comprehensive system offers a streamlined and effective solution for image-to-speech conversion, demonstrating promising results in both research and practical applications.



# TABLE OF CONTENTS

---

<b>TITLE</b>	<b>Page No</b>
1. Introduction	
1.1 Overview	1
1.2 Problem Statement	2
1.3 objectives	2
2. Literature Survey	3-7
3. Methodology	8-10
3.1 Hardware Requirements	11
3.2 Software Requirements	11
3.3 RL Algorithm	12-14
3.4 Frontend working	15
4. Sample code	16-34
4.1 RL.py	
4.2 Main.py	
4.3 Train.py	
4.4 Inference.py	
5. Results and output	35
5.1 Contribution	36
5.2 Snapshots	37
6. Conclusion	40
7. References	41

## LIST OF FIGURES

FIGURE .NO	FIGURE NAME	PAGE NO
3.1	Methodology	8
3.3.1	Reinforcement learning	14
5.2.1	Sample output 1	37
5.2.3	Sample output 2	38
5.2.4	Sample output 3	39

## CHAPTER 1

### INTRODUCTION

#### 1.1 OVERVIEW

Our project focuses on the development of a comprehensive system for image-to-speech conversion using machine learning techniques. The ability to convert visual information into spoken language has immense potential in various domains, including accessibility, education, and assistive technology. With the rapid advancement of machine learning and computer vision technologies, such a system can greatly benefit individuals with visual impairments or those who prefer auditory information presentation.

In this report, we detail our approach to building this system, which consists of two main components: image-to-text conversion and text-to-speech synthesis. For the image-to-text conversion, we leverage the power of a Generative Image-to-text Transformer (GIT) model, which streamlines the architecture into a single image encoder and text decoder. This model has been trained on a large corpus of data to accurately generate textual descriptions from input images. Moreover, we introduce an innovative reinforcement learning model alongside GIT to further enhance the performance and adaptability of our system. By incorporating reinforcement learning techniques, we aim to optimize the system's ability to generate accurate textual descriptions from images and refine its performance over time.

we extend the functionality of our system by implementing text-to-speech (TTS) synthesis, enabling the conversion of textual descriptions into spoken language. This allows for a seamless transition from visual information captured in images to auditory output, enhancing accessibility and usability for individuals with visual impairments.

Moreover, to provide real-time image-to-speech conversion, we incorporate camera access into our system, allowing users to capture images directly through the device's camera and receive spoken descriptions instantaneously. This feature enhances the practicality and versatility of our system, making it suitable for a wide range of applications, including navigation aids, object recognition, and educational tools.

## 1.2 PROBLEM STATEMENT

Despite advancements in machine learning and computer vision, existing image-to-speech conversion systems lack accuracy, real-time performance, and comprehensive functionality, limiting accessibility for visually impaired individuals. Our project aims to address this gap by developing a comprehensive system that seamlessly converts visual information into spoken language using a Generative Image-to-text Transformer(GIT) model, innovative reinforcement learning model and text-to-speech (TTS) synthesis. By incorporating camera access and reinforcement learning techniques, we aim to create an adaptable, real-time solution suitable for various applications, thereby enhancing accessibility and usability for individuals with visual impairments.

## 1.3 OBJECTIVES

- **Elevate Captioning Precision:** Strive for heightened accuracy in captions through continuous refinement and innovation within the Image Captioning Recognition (ICR) model.
- **Optimize User Experience:** Craft and fine-tune the overall user experience by prioritizing the clarity, coherence, and contextual relevance of the generated speech output to ensure an enriched and seamless interaction.
- **Foster Adaptability Across Image Varieties:** Engineer the model to exhibit robust adaptability, enabling it to effectively generalize and cater to a broad spectrum of diverse image content, ensuring versatility and reliability in various scenarios.

## CHAPTER 2

### Literature Survey

**[1] Title:** Text and Speech Recognition for Visually Impaired People using Optical Character Recognition (OCR), text-to-speech synthesizer (TTS)

**Authors:** Dr.Sujatha. K, Shruti P. Pati

#### **Abstract:**

The research by Dr. Sujatha K. and Shruti P. Pati focuses on developing a system aimed at assisting visually impaired individuals by integrating Optical Character Recognition (OCR) and text-to-speech (TTS) technologies. The system is designed to convert printed text into spoken words, thereby providing an accessible means for the visually impaired to read textual content. The study outlines the technical components involved in creating such a system, including image processing, OCR algorithms, and TTS synthesizers. The primary goal is to enhance the independence and quality of life for visually impaired users by facilitating their access to printed information through auditory means.

#### **Conclusion :**

The study concludes that the integration of OCR and TTS technologies can significantly benefit visually impaired individuals by providing them with a reliable and efficient tool for reading printed text. The developed system successfully converts textual content into speech, enabling users to access information that would otherwise be inaccessible. The authors highlight the importance of such assistive technologies in promoting inclusivity and improving the daily lives of visually impaired persons. Further research and development are encouraged to refine the system's accuracy and usability, ensuring broader application and better user experience.

**[2] Title:** Text Extraction from Image, Word, PDF and Text-to-Speech Conversion

**Authors:** Pooja Bendale<sup>1</sup>, Sanika Badhe, Sarthak Bhagat

**Abstract:**

The paper "Text Extraction from Image, Word, PDF and Text-to-Speech Conversion" by Pooja Bendale, Sanika Badhe, and Sarthak Bhagat focuses on developing a system that extracts text from images, Word documents, and PDFs, converting it into editable formats and speech. The methodology utilizes Optical Character Recognition (OCR) and text-to-speech synthesis to facilitate the conversion. This technology aims to assist users, particularly those with visual impairments, by providing an efficient means to access textual information through audio output.

**Conclusion:**

The proposed method effectively detects and extracts text from various sources, demonstrating high accuracy across different text sizes, styles, and colors. Despite its success, the system faces challenges with blurred or very small text regions. Future work will focus on enhancing the algorithm's robustness and expanding its capabilities to handle more complex and diverse text environments.

**[3] Title:** Improving Preprocessing for Image Captioning with Text-Guided Image Generation

**Authors:** Rozelle Jain, Shantanu

**Abstract:**

The paper "Image to Speech Conversion using Digital Image Processing" by Rozelle Jain and Shantanu Das discusses a system that transforms text from images into audible speech. Utilizing digital image processing techniques, the system involves capturing images, extracting text via Optical Character Recognition (OCR), and converting this text into speech using text-to-speech (TTS) synthesis. This innovation aims to aid visually impaired individuals by enabling them to access textual information audibly.

**Conclusion:**

The study concludes that the proposed system effectively bridges the gap for visually impaired users by converting visual text into speech, enhancing accessibility. While the system demonstrates good performance in controlled conditions, further improvements are needed to handle diverse real-world scenarios such as varying lighting conditions and text fonts. Future enhancements will focus on refining text extraction accuracy and expanding language support.

**[4] Title:** An Integrated Model for Text to Text, Image to Text and Audio to Text Linguistic Conversion using Machine Learning Approach

**Authors:** Aman Raj Singh, Prof. Diwakar Bhardwaj, Mridal Dixit, Lalit kumar

**Abstract:**

The paper by Aman Raj Singh, Prof. Diwakar Bhardwaj, Mridal Dixit, and Lalit Kumar introduces an integrated model designed to perform linguistic conversions across three modalities: text-to-text, image-to-text, and audio-to-text using machine learning approaches. The objective is to create a versatile system that can process and convert different types of input data into text, thereby enhancing accessibility and facilitating efficient information processing. The model leverages advanced machine learning techniques, including natural language processing (NLP) for text-to-text conversions, Optical Character Recognition (OCR) for converting images to text, and automated speech recognition (ASR) for converting audio inputs to text. The study provides a detailed explanation of the system's architecture, the algorithms used, and the integration process that allows seamless conversion across different data types.

**Conclusion :**

The study concludes that the proposed integrated model achieves high performance in converting text, image, and audio inputs into text, demonstrating the effectiveness of using machine learning techniques for linguistic conversions. The results indicate that the system is capable of accurately processing and converting diverse forms of data, which has significant implications for enhancing accessibility, especially for individuals with disabilities, and for improving various applications such as data entry automation and human-computer interaction. The authors emphasize the importance of further refining the model to enhance its accuracy, scalability, and user-friendliness. Future research directions include optimizing the machine learning algorithms used, improving the integration of different modalities, and expanding the system's capabilities to handle more complex linguistic conversion tasks.



**[5] Title:** OCR Based Image Text to Speech Conversion using K-Nearest Neighbors and Comparing with Fuzzy K-Means Clustering Algorithm.

**Authors:** M.Pandu Babu, Anitha

**Abstract:**

The paper by M. Pandu Babu and Anitha G. presents a study on converting text from images to speech using Optical Character Recognition (OCR) enhanced by machine learning techniques. The primary objective is to develop a system that accurately reads and vocalizes text found in images, thereby aiding visually impaired individuals and automating data entry processes. The research employs the K-Nearest Neighbors (KNN) algorithm for the OCR process and compares its performance with the Fuzzy K-Means Clustering algorithm. The study aims to evaluate the effectiveness and efficiency of these algorithms in terms of accuracy and processing time. Detailed experimental results and analyses are provided to highlight the strengths and weaknesses of each approach.

**Conclusion:**

The study concludes that both the K-Nearest Neighbors and Fuzzy K-Means Clustering algorithms are viable for OCR-based image text to speech conversion, but they exhibit different strengths. The KNN algorithm demonstrates higher accuracy in recognizing and converting text, making it more suitable for applications where precision is critical. On the other hand, the Fuzzy K-Means Clustering algorithm shows advantages in terms of processing speed and computational efficiency. The comparison underscores the importance of selecting the appropriate algorithm based on specific application requirements. The authors suggest that further research could explore hybrid approaches that combine the strengths of both algorithms to achieve better overall performance in OCR-based text to speech systems.

## CHAPTER 3

### METHODOLOGY

The methodology of the provided code involves several components for building an image-to-speech conversion application using a deep learning model. Here's a breakdown of the methodology based on the code provided:

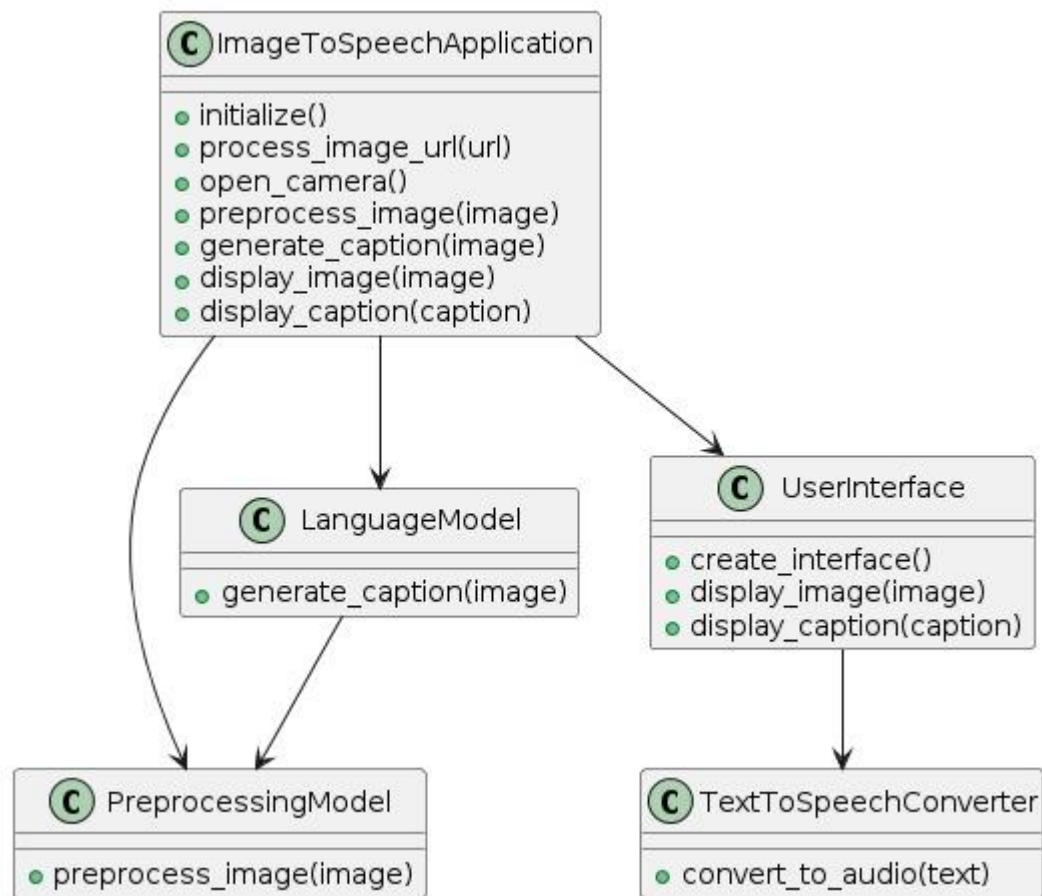


Fig 3.1 Methodology

### **1. Image Processing and Model Initialization:**

- The application initializes a processor and a model for image captioning from the Hugging Face Transformers library.
- The processor is responsible for preprocessing images before feeding them into the model.
- The model is a pretrained language model capable of generating captions for images.

### **2. User Interface Setup:**

- The application uses Tkinter, a Python GUI toolkit, to create a user-friendly interface.
- The interface includes an entry field for entering image URLs, buttons for processing URLs or opening the camera, a text area for displaying speech output, and a label for displaying images.

### **3. Processing Image URLs:**

- When the user enters an image URL and clicks the "Process URL" button, the application retrieves the image from the URL.
- The retrieved image is then displayed in the interface, and its corresponding speech output is generated using the pretrained model.
- The generated speech is displayed in the text area, and text-to-speech conversion is performed to audibly output the speech.

### **4. Opening Camera:**

- Clicking the "Open Camera" button activates the device's camera using OpenCV.
- The camera feed is displayed in real-time in the interface.
- The application continuously captures frames from the camera, processes them, generates speech output, and updates the interface accordingly.

### **5. Generating Speech Output:**

- For both image URLs and camera feed, the captured image is resized to a fixed size (224x224 pixels) to match the input size expected by the pretrained model.
- The pixel values of the resized image are preprocessed using the initialized processor.
- The processed image is then fed into the pretrained model to generate a caption for the image.
- The generated caption is displayed in the text area and converted into speech using the pyttsx3 library for audible output.

### **6. Exception Handling:**

- The application includes error handling mechanisms to catch and display any exceptions that may occur during image processing, model inference, or camera initialization.
- Error messages are displayed using message boxes to inform the user about the nature of the error.

The methodology involves integrating image processing, deep learning model inference, user interface design, and error handling to create an interactive image-to-speech conversion application. The application allows users to input image URLs or use their device's camera to capture images, which are then processed to generate descriptive speech output.

## **SYSTEM REQUIREMENTS**

### **3.1 Hardware Requirements:**

CPU: Intel CORE i5 8<sup>th</sup> Gen (Clock speed: 2.8GHz)

Hard Disk: 120GB

RAM : 8 GB DDR4

### **3.2 Software Requirements:**

Operating System: Windows 10,11

Scripting Language: Python 3.11.4 64 bit

Deep learning Library: Tensorflow /Keras

Front End: Python/tkinter

Text Editor: Visual Studio Code

### 3.3 RL – Algorithm :

#### 1. Initialization:

- Necessary libraries are imported, including TensorFlow and OpenAI Gym.
- An environment class Image To TextEnv is defined, which simulates the task of converting an image to text.
- A neural network model class RLModel is defined, which will be used to approximate the agent's policy.

#### 2. Environment:

- The Image To TextEnv class defines the environment in which the agent operates.
- It has methods for resetting the environment (reset), taking actions (step), and rendering the environment (render).

#### 3. RL Model:

- The RLModel class defines the neural network architecture that serves as the policy network for the agent.
- It consists of a combination of layers, including a flattening layer, dense layers, and an LSTM layer.
- The model outputs action probabilities based on the current state.

#### 4. Training Loop:

- Hyperparameters such as learning rate, discount factor (gamma), number of episodes, and exploration rate (epsilon) are defined.
- The environment is instantiated with an API key.
- The RL model is instantiated.
- The training loop runs for a fixed number of episodes.
- For each episode:
- The environment is reset.
- The agent interacts with the environment for a maximum number of timesteps.

- At each timestep, the agent selects an action according to its policy (with  $\epsilon$ -greedy exploration).
- The action is taken in the environment, and the next state, reward, and termination signal are received.
- The discounted reward is computed.
- The REINFORCE loss is calculated based on the log probability of the selected action and the discounted reward.
- The gradients of the loss with respect to the model's parameters are computed using automatic differentiation.
- The model parameters are updated using the Adam optimizer.
- The loop continues until the episode terminates.

### 5. Model Saving:

- After training, the trained RL model is saved to disk.

The REINFORCE algorithm is a policy gradient method that aims to maximize the expected cumulative reward by directly optimizing the policy. It does so by estimating the gradient of the expected reward with respect to the policy parameters and updating the parameters in the direction that increases the expected reward. In this implementation, the REINFORCE loss is computed as the negative log probability of the selected action multiplied by the discounted reward.

The **total cumulative reward** of the agent. The figure below illustrates the **action-reward feedback loop** of a generic RL model.

Some key terms that describe the basic elements of an RL problem are:

1. **Environment** — Physical world in which the agent operates
2. **State** — Current situation of the agent
3. **Reward** — Feedback from the environment
4. **Policy** — Method to map agent's state to actions
5. **Value** — Future reward that an agent would receive by taking an action in a particular state

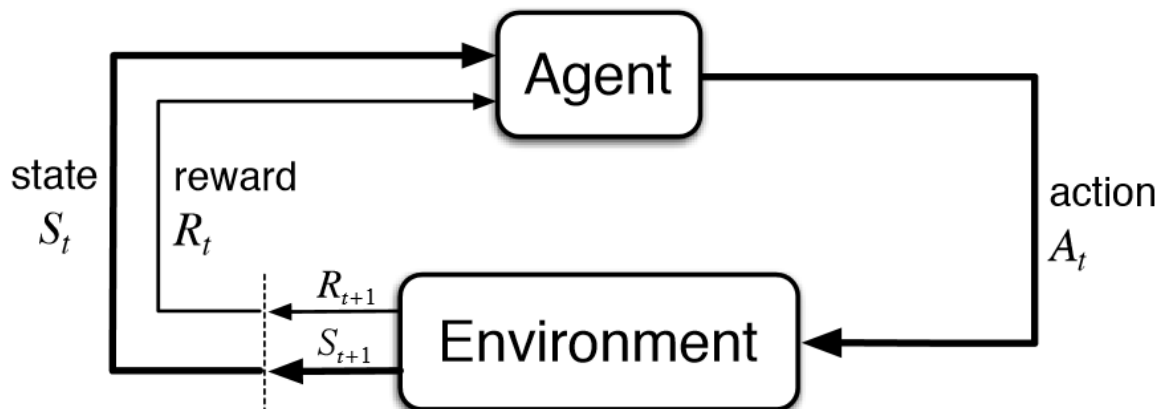


Fig 3.3.1 Reinforcement Learning



### 3.4 FRONTEND WORKING

**1. Initialization:** The ImageToSpeechApp class initializes the Tkinter root window (root) and sets its title to "Image to Speech Converter".

**2. URL Entry and Buttons:** There's an entry widget (self.url\_entry) where users can input the URL of an image. Two buttons are provided: "Process URL" (self.url\_button) and "Open Camera" (self.camera\_button). When the "Process URL" button is clicked, it calls the process\_url method, and when the "Open Camera" button is clicked, it calls the open\_camera method.

**3. Speech Output and Image Display:** A label (self.speech\_label) displays "Speech Output". A text widget (self.speech\_text) displays the speech generated from the image. An empty label (self.image\_label) is initially present, which will display the processed image.

**4. Model Initialization:** The Hugging Face transformers library is used to initialize the processor and model for generating captions from images. The processor (self.processor) is initialized to preprocess the image, and the model (self.model) is initialized to generate captions.

**5. Processing URL and Camera Feed:** When the "Process URL" button is clicked (process\_url method), the URL entered is fetched and an attempt is made to download the image from that URL. If successful, the image is displayed, and its caption is generated and spoken. When the "Open Camera" button is clicked (open\_camera method), it initializes the camera, continuously captures frames, displays them, generates captions for the frames, and speaks them out.

**6. Image Display:** The display\_image method is responsible for resizing the image to fit within a specific size and displaying it in the GUI using a label widget.

**7. Speech Generation:** The generate\_speech\_from\_image method processes the image (resizing it appropriately), generates a caption using the model, updates the speech output text widget with the generated caption, and converts the text to speech using the pyttsx3 library.

**8. Exception Handling:** Exceptions are handled throughout the code using try-except blocks. If any error occurs during image processing or speech generation, an error message is displayed using a message box.

## CHAPTER 4

### SAMPLE CODE

#### 4.1 REINFORCEMENT LEARNING MODEL

```
import numpy as np
import tensorflow as tf
from tensorflow.keras.layers import Dense, LSTM, Flatten
from tensorflow.keras.optimizers import Adam
import gym

# Define RL environment
class ImageToTextEnv(gym.Env):
    def __init__(self, api_key):
        super(ImageToTextEnv, self).__init__()
        self.api_key = api_key
        self.action_space = gym.spaces.Discrete(10)
        self.observation_space = gym.spaces.Box(low=0, high=255, shape=(224, 224,
3), dtype=np.uint8)
        # Other initialization, such as API setup

    def reset(self):
        # Reset environment
        return np.zeros((224, 224, 3), dtype=np.uint8)

    def step(self, action):
        reward = 1 if action == 0 else -1
        next_state = np.zeros((224, 224, 3), dtype=np.uint8)
        done = False
        info = {}
        return next_state, reward, done, info

    def render(self):
        # Optional: visualize environment state
        pass
```

```
# Define RL model
class RLModel(tf.keras.Model):
    def __init__(self, action_space, name='rl_model'):
        super(RLModel, self).__init__(name=name)
        self.action_space = action_space
        self.flatten = Flatten()
        self.dense1 = Dense(64, activation='relu')
        self.lstm = LSTM(128)
        self.dense2 = Dense(action_space, activation='softmax')

    def call(self, inputs):
        x = self.flatten(inputs)
        x = self.dense1(x)
        x = self.lstm(x[:, tf.newaxis, :])
        return self.dense2(x)

    def get_config(self):
        config = super(RLModel, self).get_config()
        config.update({'action_space': self.action_space})
        return config

    @classmethod
    def from_config(cls, config):
        name = config.pop('name', 'rl_model')
        action_space = config.pop('action_space')
        return cls(action_space, name=name)

# Hyperparameters
learning_rate = 0.001
gamma = 0.99
num_episodes = 50
epsilon = 0.1 # Epsilon for  $\epsilon$ -greedy exploration

# Create environment
api_key = "K88153881288957"
env = ImageToTextEnv(api_key)

# Create RL model
model = RLModel(env.action_space.n) # Here we directly pass the number of actions
optimizer = Adam(learning_rate)

# Training loop
for episode in range(num_episodes):
    state = env.reset()
```

```
episode_reward = 0

for timestep in range(50):
    with tf.GradientTape() as tape:
        # Convert state to suitable format
        state = np.array(state)[np.newaxis, :]

        # Get action probabilities from model
        action_probs = model(state)
        if np.random.rand() < epsilon:
            action = np.random.randint(env.action_space.n)
        else:
            action = np.argmax(action_probs.numpy()[0])

        # Take action, get next state and reward from environment
        next_state, reward, done, _ = env.step(action)

        # Compute discounted reward
        episode_reward += reward
        discounted_reward = reward * (gamma ** timestep)

        # Apply REINFORCE loss
        loss = -tf.math.log(action_probs[0, action]) * discounted_reward

    # Backpropagation
    gradients = tape.gradient(loss, model.trainable_variables)
    optimizer.apply_gradients(zip(gradients, model.trainable_variables))

    if done:
        break
    state = next_state

# Print episode info
print(f"Episode {episode + 1}: Total Reward = {episode_reward}")

# Save trained model
model.save("image_to_text_rl_model.h5")
```

## 4.2 MAIN.PY CODE

```
import numpy as np
import tensorflow as tf
from tensorflow.keras.models import load_model
from tensorflow.keras.layers import Dense, LSTM, Flatten
from tensorflow.keras.optimizers import Adam
import gym
import cv2

# Define RL environment
class ImageToTextEnv(gym.Env):
    def __init__(self, api_key):
        super(ImageToTextEnv, self).__init__()
        self.api_key = api_key
        self.action_space = gym.spaces.Discrete(10)
        self.observation_space = gym.spaces.Box(low=0, high=255, shape=(224, 224, 3), dtype=np.uint8)
        # Other initialization, such as API setup

    def reset(self):
        # Reset environment
        return np.zeros((224, 224, 3), dtype=np.uint8)

    def step(self, action):
        reward = 1 if action == 0 else -1
        next_state = np.zeros((224, 224, 3), dtype=np.uint8)
        done = False
        info = {}
        return next_state, reward, done, info

    def render(self):
        # Optional: visualize environment state
        pass
```

```
# Define RL model
class RLModel(tf.keras.Model):
    def __init__(self, action_space):
        super(RLModel, self).__init__()
        self.flatten = Flatten()
        self.dense1 = Dense(64, activation='relu')
        self.lstm = LSTM(128)
        self.dense2 = Dense(action_space, activation='softmax')

    def call(self, inputs):
        x = self.flatten(inputs)
        x = self.dense1(x)
        x = self.lstm(x[:, tf.newaxis, :]) # Reshape to (batch_size, timesteps,
input_dim)
        return self.dense2(x)

# Define custom objects to load the model
custom_objects = {'RLModel': RLModel}

# Load the model with custom objects
model = load_model("image_to_text_rl_model.h5", custom_objects=custom_objects)

# Initialize OpenCV video capture
cap = cv2.VideoCapture(0)

# Main loop for real-time captioning
while True:
    ret, frame = cap.read() # Read frame from webcam
    if not ret:
        break
    # Preprocess frame (replace this with your actual preprocessing function)
    preprocessed_frame = frame

    # Generate caption
    caption = model.predict(np.expand_dims(preprocessed_frame, axis=0))
    print("Caption:", caption) # Print caption (you might want to display it
differently)
    # Display the frame with the caption
    cv2.imshow('Frame', frame)

    # Exit if 'q' is pressed
    if cv2.waitKey(1) & 0xFF == ord('q'):
        break

# Release video capture and close OpenCV windows
cap.release()
cv2.destroyAllWindows()
```

## 4.3 Train.py

```
import json
import os.path as op
from .common import qd_tqdm as tqdm
from .common import json_dump
from .common import pilimg_from_base64
from .common import get_mpi_rank, get_mpi_size, get_mpi_local_rank

from .tsv_io import TSVFile, tsv_writer, tsv_reader
from .common import write_to_file
import torch
import PIL
from pprint import pformat
import logging
from transformers import BertTokenizer
from torchvision.transforms import Compose, Resize, CenterCrop, ToTensor,
Normalize
from PIL import Image
from azfuse import File

from .common import init_logging
from .common import parse_general_args
from .tsv_io import load_from_yaml_file
from .torch_common import torch_load
from .torch_common import load_state_dict
from .process_image import load_image_by_pil
from .model import get_git_model

class MinMaxResizeForTest(object):
    def __init__(self, min_size, max_size):
        self.min_size = min_size
        self.max_size = max_size

    def get_size(self, image_size):
        w, h = image_size
        size = self.min_size
        max_size = self.max_size

        min_original_size = float(min((w, h)))
        max_original_size = float(max((w, h)))
        if max_original_size / min_original_size * size > max_size:
            size = int(round(max_size * min_original_size / max_original_size))

        if (w <= h and w == size) or (h <= w and h == size):
            return (h, w)
```

```
        if w < h:
            ow = size
            oh = int(size * h / w)
        else:
            oh = size
            ow = int(size * w / h)

    return (oh, ow)

def _repr_(self):
    return 'MinMaxResizeForTest({}, {})'.format(
        self.min_size, self.max_size)

def _call_(self, img):
    size = self.get_size(img.size)
    import torchvision.transforms.functional as F
    image = F.resize(img, size, interpolation=PIL.Image.BICUBIC)
    return image

def test_git_inference_single_image(image_path, model_name, prefix):
    param = {}
    if File.isfile(f'aux_data/models/{model_name}/parameter.yaml'):
        param =
load_from_yaml_file(f'aux_data/models/{model_name}/parameter.yaml')

    tokenizer = BertTokenizer.from_pretrained('bert-base-uncased',
do_lower_case=True)
    if isinstance(image_path, str):
        image_path = [image_path]
    # if it is more than 1 image, it is normally a video with multiple image
    # frames
    img = [load_image_by_pil(i) for i in image_path]

    transforms = get_image_transform(param)
    img = [transforms(i) for i in img]

    # model
    model = get_git_model(tokenizer, param)
    pretrained = f'output/{model_name}/snapshot/model.pt'
    checkpoint = torch_load(pretrained)['model']
    load_state_dict(model, checkpoint)
    model.cuda()
    model.eval()
    img = [i.unsqueeze(0).cuda() for i in img]

    # prefix
```



```

max_text_len = 40
prefix_encoding = tokenizer(prefix,
                             padding='do_not_pad',
                             truncation=True,
                             add_special_tokens=False,
                             max_length=max_text_len)
payload = prefix_encoding['input_ids']
if len(payload) > max_text_len - 2:
    payload = payload[-(max_text_len - 2):]
input_ids = [tokenizer.cls_token_id] + payload

with torch.no_grad():
    result = model({
        'image': img,
        'prefix': torch.tensor(input_ids).unsqueeze(0).cuda(),
    })
    cap = tokenizer.decode(result['predictions'][0].tolist(),
skip_special_tokens=True)
    logging.info('output: {}'.format(cap))

def get_image_transform(param):
    crop_size = param.get('test_crop_size', 224)
    if 'test_respect_ratio_max' in param:
        trans = [
            MinMaxResizeForTest(crop_size, param['test_respect_ratio_max'])
        ]
    else:
        trans = [
            Resize(crop_size, interpolation=Image.BICUBIC),
            CenterCrop(crop_size),
            lambda image: image.convert("RGB"),
        ]
    trans.extend([
        ToTensor(),
        Normalize(
            (0.48145466, 0.4578275, 0.40821073),
            (0.26862954, 0.26130258, 0.27577711),
        ),
    ])
    transforms = Compose(trans)
    return transforms

def test_git_inference_single_tsv(image_tsv, model_name, question_tsv, out_tsv):
    param = {}
    if File.isfile(f'output/{model_name}/parameter.yaml'):
        param = load_from_yaml_file(f'output/{model_name}/parameter.yaml')

```

```
tokenizer = BertTokenizer.from_pretrained('bert-base-uncased',
do_lower_case=True)
image_tsv = TSVFile(image_tsv)
question_tsv = TSVFile(question_tsv) if question_tsv else None

transforms = get_image_transform(param)

# model
model = get_git_model(tokenizer, param)
pretrained = f'output/{model_name}/snapshot/model.pt'
checkpoint = torch_load(pretrained)['model']
load_state_dict(model, checkpoint)
model.eval()

torch.cuda.set_device(get_mpi_local_rank())
model.cuda()

# prefix
max_text_len = 40
rank = get_mpi_rank()
world_size = get_mpi_size()
def get_rank_specific_tsv(rank):
    return '{}.{}.{}.tsv'.format(out_tsv, rank, world_size)
if world_size > 1:
    curr_out_tsv = get_rank_specific_tsv(rank)
else:
    curr_out_tsv = out_tsv
total = len(image_tsv)
curr_size = (total + world_size - 1) // world_size
curr_start = curr_size * rank
curr_end = curr_start + curr_size
curr_end = min(curr_end, total)

if question_tsv:
    def gen_rows():
        for i in tqdm(range(curr_start, curr_end)):
            image_key, image_col = image_tsv[i]
            q_key, q_info = question_tsv[i]
            assert image_key == q_key
            q_info = json.loads(q_info)
            img = pilimg_from_base64(image_col)
            img = transforms(img)
            img = img.cuda().unsqueeze(0)
            for q in q_info:
                prefix = q['question']
                prefix_encoding = tokenizer(prefix,
                                                padding='do_not_pad',
                                                truncation=True,
```

```

                                add_special_tokens=False,
                                max_length=max_text_len)
    payload = prefix_encoding['input_ids']
    if len(payload) > max_text_len - 2:
        payload = payload[-(max_text_len - 2):]
    input_ids = [tokenizer.cls_token_id] + payload
    with torch.no_grad():
        result = model({
            'image': img,
            'prefix': torch.tensor(input_ids).unsqueeze(0).cuda(),
        })
    answer = tokenizer.decode(result['predictions'][0].tolist(),
skip_special_tokens=True)
    result = {'answer': answer, 'question_id': q['question_id']}
    yield json_dump(result),
else:
    def gen_rows():
        for i in tqdm(range(curr_start, curr_end)):
            key, col = image_tsv[i]
            img = pilimg_from_base64(col)
            img = transforms(img)
            img = img.cuda().unsqueeze(0)
            with torch.no_grad():
                result = model({
                    'image': img,
                })
            cap = tokenizer.decode(result['predictions'][0].tolist(),
skip_special_tokens=True)
            yield key, json_dump([{'caption': cap}])
    tsv_writer(gen_rows(), curr_out_tsv)
    if world_size > 1 and rank == 0:
        all_sub_tsv = [get_rank_specific_tsv(i) for i in range(world_size)]
        while True:
            not_ready = [t for t in all_sub_tsv if not File.isfile(t)]
            if len(not_ready) == 0:
                break
            else:
                import time
                logging.info('waiting {}'.format(','.join(not_ready)))
                time.sleep(5)
        from .tsv_io import concat_tsv_files
        concat_tsv_files(all_sub_tsv, out_tsv)

def convert_tsv_to_vqa_json(predict_file, out_json):
    result = [json.loads(s) for s, in tsv_reader(predict_file)]
    write_to_file(json_dump(result), out_json)

def convert_tsv_to_coco_format(res_tsv, outfile,

```

```

        sep='\t', key_col=0, cap_col=1):
    results = []
    with open(res_tsv) as fp:
        for line in fp:
            parts = line.strip().split(sep)
            key = parts[key_col]
            if cap_col < len(parts):
                caps = json.loads(parts[cap_col])
                if len(caps) == 0:
                    caps = [{'caption': ''}]
                assert len(caps) == 1, 'cannot evaluate multiple captions per
image'
                cap = caps[0]['caption']
            else:
                # empty caption generated
                cap = ""
            results.append(
                {'image_id': key,
                 'caption': cap}
            )
    with open(outfile, 'w') as fp:
        json.dump(results, fp)

def iter_caption_to_json(iter_caption, json_file):
    # save gt caption to json format so that we can call the api
    key_captions = [(key, json.loads(p)) for key, p in iter_caption]

    info = {
        'info': 'dummy',
        'licenses': 'dummy',
        'type': 'captions',
    }
    info['images'] = [{'file_name': k, 'id': k} for k, _ in key_captions]
    n = 0
    annotations = []
    for k, cs in key_captions:
        for c in cs:
            annotations.append({
                'image_id': k,
                'caption': c['caption'],
                'id': n
            })
            n += 1
    info['annotations'] = annotations
    write_to_file(json.dumps(info), json_file)

def evaluate_on_coco_caption(res_file, label_file, outfile=None,
):

```

```
if not outfile:
    outfile = op.splitext(res_file)[0] + '.eval.json'

if res_file.endswith('.tsv'):
    res_file_coco = op.splitext(res_file)[0] + '_coco_format.json'
    convert_tsv_to_coco_format(res_file, res_file_coco)
else:
    res_file_coco = res_file

if label_file.endswith('.tsv'):
    json_caption = '/tmp/{}'.format(label_file)
    iter_caption_to_json(
        TSVFile(label_file),
        json_caption)
    label_file = json_caption

from pycocotools.coco import COCO
from pycocoevalcap.eval import COCOEvalCap

coco = COCO(label_file)
cocoRes = coco.loadRes(res_file_coco)
cocoEval = COCOEvalCap(coco, cocoRes)
# evaluate results
# SPICE will take a few minutes the first time, but speeds up due to caching
cocoEval.evaluate()
result = cocoEval.eval
if not outfile:
    print(result)
else:
    with open(outfile, 'w') as fp:
        json.dump(result, fp, indent=4)
return result

if __name__ == '__main__':
    init_logging()
    kwargs = parse_general_args()
    logging.info('param:\n{}'.format(pformat(kwargs)))
    function_name = kwargs['type']
    del kwargs['type']
    locals()[function_name](**kwargs)
```

## 4.4 INFERENCE.PY

```
import json
import os.path as op
from .common import qd_tqdm as tqdm
from .common import json_dump
from .common import pilimg_from_base64
from .common import get_mpi_rank, get_mpi_size, get_mpi_local_rank

from .tsv_io import TSVFile, tsv_writer, tsv_reader
from .common import write_to_file
import torch
import PIL
from pprint import pformat
import logging
from transformers import BertTokenizer
from torchvision.transforms import Compose, Resize, CenterCrop, ToTensor,
Normalize
from PIL import Image
from azfuse import File

from .common import init_logging
from .common import parse_general_args
from .tsv_io import load_from_yaml_file
from .torch_common import torch_load
from .torch_common import load_state_dict
from .process_image import load_image_by_pil
from .model import get_git_model

class MinMaxResizeForTest(object):
    def __init__(self, min_size, max_size):
        self.min_size = min_size
        self.max_size = max_size

    def get_size(self, image_size):
        w, h = image_size
        size = self.min_size
        max_size = self.max_size

        min_original_size = float(min((w, h)))
        max_original_size = float(max((w, h)))
        if max_original_size / min_original_size * size > max_size:
            size = int(round(max_size * min_original_size / max_original_size))

        if (w <= h and w == size) or (h <= w and h == size):
            return (h, w)
```

```
        if w < h:
            ow = size
            oh = int(size * h / w)
        else:
            oh = size
            ow = int(size * w / h)

    return (oh, ow)

def _repr_(self):
    return 'MinMaxResizeForTest({}, {})'.format(
        self.min_size, self.max_size)

def _call_(self, img):
    size = self.get_size(img.size)
    import torchvision.transforms.functional as F
    image = F.resize(img, size, interpolation=PIL.Image.BICUBIC)
    return image

def test_git_inference_single_image(image_path, model_name, prefix):
    param = {}
    if File.isfile(f'aux_data/models/{model_name}/parameter.yaml'):
        param =
load_from_yaml_file(f'aux_data/models/{model_name}/parameter.yaml')

    tokenizer = BertTokenizer.from_pretrained('bert-base-uncased',
do_lower_case=True)
    if isinstance(image_path, str):
        image_path = [image_path]
    # if it is more than 1 image, it is normally a video with multiple image
    # frames
    img = [load_image_by_pil(i) for i in image_path]

    transforms = get_image_transform(param)
    img = [transforms(i) for i in img]

    # model
    model = get_git_model(tokenizer, param)
    pretrained = f'output/{model_name}/snapshot/model.pt'
    checkpoint = torch_load(pretrained)['model']
    load_state_dict(model, checkpoint)
    model.cuda()
    model.eval()
    img = [i.unsqueeze(0).cuda() for i in img]

    # prefix
```

```

max_text_len = 40
prefix_encoding = tokenizer(prefix,
                             padding='do_not_pad',
                             truncation=True,
                             add_special_tokens=False,
                             max_length=max_text_len)
payload = prefix_encoding['input_ids']
if len(payload) > max_text_len - 2:
    payload = payload[-(max_text_len - 2):]
input_ids = [tokenizer.cls_token_id] + payload

with torch.no_grad():
    result = model({
        'image': img,
        'prefix': torch.tensor(input_ids).unsqueeze(0).cuda(),
    })
    cap = tokenizer.decode(result['predictions'][0].tolist(),
skip_special_tokens=True)
    logging.info('output: {}'.format(cap))

def get_image_transform(param):
    crop_size = param.get('test_crop_size', 224)
    if 'test_respect_ratio_max' in param:
        trans = [
            MinMaxResizeForTest(crop_size, param['test_respect_ratio_max'])
        ]
    else:
        trans = [
            Resize(crop_size, interpolation=Image.BICUBIC),
            CenterCrop(crop_size),
            lambda image: image.convert("RGB"),
        ]
    trans.extend([
        ToTensor(),
        Normalize(
            (0.48145466, 0.4578275, 0.40821073),
            (0.26862954, 0.26130258, 0.27577711),
        ),
    ])
    transforms = Compose(trans)
    return transforms

def test_git_inference_single_tsv(image_tsv, model_name, question_tsv, out_tsv):
    param = {}
    if File.isfile(f'output/{model_name}/parameter.yaml'):
        param = load_from_yaml_file(f'output/{model_name}/parameter.yaml')

```



```

tokenizer = BertTokenizer.from_pretrained('bert-base-uncased',
do_lower_case=True)
image_tsv = TSVFile(image_tsv)
question_tsv = TSVFile(question_tsv) if question_tsv else None

transforms = get_image_transform(param)

# model
model = get_git_model(tokenizer, param)
pretrained = f'output/{model_name}/snapshot/model.pt'
checkpoint = torch_load(pretrained)['model']
load_state_dict(model, checkpoint)
model.eval()

torch.cuda.set_device(get_mpi_local_rank())
model.cuda()

# prefix
max_text_len = 40
rank = get_mpi_rank()
world_size = get_mpi_size()
def get_rank_specific_tsv(rank):
    return '{}.{}.{}.tsv'.format(out_tsv, rank, world_size)
if world_size > 1:
    curr_out_tsv = get_rank_specific_tsv(rank)
else:
    curr_out_tsv = out_tsv
total = len(image_tsv)
curr_size = (total + world_size - 1) // world_size
curr_start = curr_size * rank
curr_end = curr_start + curr_size
curr_end = min(curr_end, total)

if question_tsv:
    def gen_rows():
        for i in tqdm(range(curr_start, curr_end)):
            image_key, image_col = image_tsv[i]
            q_key, q_info = question_tsv[i]
            assert image_key == q_key
            q_info = json.loads(q_info)
            img = pilimg_from_base64(image_col)
            img = transforms(img)
            img = img.cuda().unsqueeze(0)
            for q in q_info:
                prefix = q['question']
                prefix_encoding = tokenizer(prefix,
                                                padding='do_not_pad',
                                                truncation=True,

```

```

                                add_special_tokens=False,
                                max_length=max_text_len)
    payload = prefix_encoding['input_ids']
    if len(payload) > max_text_len - 2:
        payload = payload[-(max_text_len - 2):]
    input_ids = [tokenizer.cls_token_id] + payload
    with torch.no_grad():
        result = model({
            'image': img,
            'prefix': torch.tensor(input_ids).unsqueeze(0).cuda(),
        })
    answer = tokenizer.decode(result['predictions'][0].tolist(),
skip_special_tokens=True)
    result = {'answer': answer, 'question_id': q['question_id']}
    yield json_dump(result),
else:
    def gen_rows():
        for i in tqdm(range(curr_start, curr_end)):
            key, col = image_tsv[i]
            img = pilimg_from_base64(col)
            img = transforms(img)
            img = img.cuda().unsqueeze(0)
            with torch.no_grad():
                result = model({
                    'image': img,
                })
            cap = tokenizer.decode(result['predictions'][0].tolist(),
skip_special_tokens=True)
            yield key, json_dump([{'caption': cap}])
    tsv_writer(gen_rows(), curr_out_tsv)
    if world_size > 1 and rank == 0:
        all_sub_tsv = [get_rank_specific_tsv(i) for i in range(world_size)]
        while True:
            not_ready = [t for t in all_sub_tsv if not File.isfile(t)]
            if len(not_ready) == 0:
                break
            else:
                import time
                logging.info('waiting {}'.format(','.join(not_ready)))
                time.sleep(5)
        from .tsv_io import concat_tsv_files
        concat_tsv_files(all_sub_tsv, out_tsv)

def convert_tsv_to_vqa_json(predict_file, out_json):
    result = [json.loads(s) for s, in tsv_reader(predict_file)]
    write_to_file(json_dump(result), out_json)

def convert_tsv_to_coco_format(res_tsv, outfile,

```

```

        sep='\t', key_col=0, cap_col=1):
    results = []
    with open(res_tsv) as fp:
        for line in fp:
            parts = line.strip().split(sep)
            key = parts[key_col]
            if cap_col < len(parts):
                caps = json.loads(parts[cap_col])
                if len(caps) == 0:
                    caps = [{'caption': ''}]
                assert len(caps) == 1, 'cannot evaluate multiple captions per
image'

                cap = caps[0]['caption']
            else:
                # empty caption generated
                cap = ""
            results.append(
                {'image_id': key,
                 'caption': cap}
            )
    with open(outfile, 'w') as fp:
        json.dump(results, fp)

def iter_caption_to_json(iter_caption, json_file):
    # save gt caption to json format so that we can call the api
    key_captions = [(key, json.loads(p)) for key, p in iter_caption]

    info = {
        'info': 'dummy',
        'licenses': 'dummy',
        'type': 'captions',
    }
    info['images'] = [{'file_name': k, 'id': k} for k, _ in key_captions]
    n = 0
    annotations = []
    for k, cs in key_captions:
        for c in cs:
            annotations.append({
                'image_id': k,
                'caption': c['caption'],
                'id': n
            })
            n += 1
    info['annotations'] = annotations
    write_to_file(json.dumps(info), json_file)

def evaluate_on_coco_caption(res_file, label_file, outfile=None,
):

```

```
if not outfile:
    outfile = op.splitext(res_file)[0] + '.eval.json'

if res_file.endswith('.tsv'):
    res_file_coco = op.splitext(res_file)[0] + '_coco_format.json'
    convert_tsv_to_coco_format(res_file, res_file_coco)
else:
    res_file_coco = res_file

if label_file.endswith('.tsv'):
    json_caption = '/tmp/{}'.format(label_file)
    iter_caption_to_json(
        TSVFile(label_file),
        json_caption)
    label_file = json_caption

from pycocotools.coco import COCO
from pycocoevalcap.eval import COCOEvalCap

coco = COCO(label_file)
cocoRes = coco.loadRes(res_file_coco)
cocoEval = COCOEvalCap(coco, cocoRes)
# evaluate results
# SPICE will take a few minutes the first time, but speeds up due to caching
cocoEval.evaluate()
result = cocoEval.eval
if not outfile:
    print(result)
else:
    with open(outfile, 'w') as fp:
        json.dump(result, fp, indent=4)
return result

if __name__ == '__main__':
    init_logging()
    kwargs = parse_general_args()
    logging.info('param:\n{}'.format(pformat(kwargs)))
    function_name = kwargs['type']
    del kwargs['type']
    locals()[function_name](**kwargs)
```

## CHAPTER 5

### RESULTS AND OUTPUT

The "Image to Speech" frontend interface provides users with a seamless and intuitive platform to convert visual information into spoken language. Upon accessing the interface, users are greeted with the option to input an image URL or utilize the camera functionality to capture real-time images. This straightforward design empowers users to choose their preferred input method, catering to a wide range of preferences and accessibility needs.

Once an image URL is entered, users can opt to process the URL or directly access the device's camera by selecting the corresponding buttons. This user-friendly interface facilitates smooth navigation and enhances the overall user experience. Users can seamlessly transition between input methods, ensuring flexibility and convenience.

Upon processing the URL or accessing the camera, the system promptly generates a spoken description of the visual content depicted in the image. This spoken output is accompanied by a textual description displayed on the interface, providing users with multiple modalities for accessing information. Additionally, the image itself is prominently featured on the interface, serving as a visual reference point and enhancing comprehension.

The integration of speech output, textual display, and image presentation creates a cohesive and multi-modal experience for users. This comprehensive approach caters to individuals with diverse preferences and accessibility needs, ensuring inclusivity and usability. By offering a seamless interface with streamlined functionalities, the "Image to Speech" frontend facilitates efficient conversion of visual information into spoken language, ultimately enhancing user engagement and satisfaction.

## 5.1 OUR MAIN CONTRIBUTIONS

**Attention Mechanism and Transformer Layers:** Our contribution lies in the implementation and optimization of attention mechanisms and transformer layers within the image captioning system. These components significantly enhance the model's ability to capture relevant features in input images and generate accurate captions, thereby improving overall performance and quality of the system.

**Reinforcement Learning (RL) Model Agent:** We contribute by integrating a reinforcement learning (RL) model agent into the image captioning system. This RL agent continuously refines and optimizes the captioning process through interaction with the environment and feedback mechanisms. By leveraging RL techniques, our system adapts and improves over time, resulting in more accurate and contextually relevant captions.

**User-Friendly Interface:** Our contribution extends to the development of a user-friendly interface that ensures ease of use and accessibility for all users. Through intuitive controls, clear visual feedback, and seamless interaction, our interface enhances the overall user experience, making the image captioning system more accessible and efficient for a wide range of users.

**Single and Multiple Image Captioning:** We contribute to the versatility of the system by enabling both single and multiple image captioning capabilities. Users can input images either through URL links or camera access, allowing for flexible usage scenarios. This functionality expands the utility of the system, catering to various needs and preferences of users.

## 5.2 SNAPSHOTS

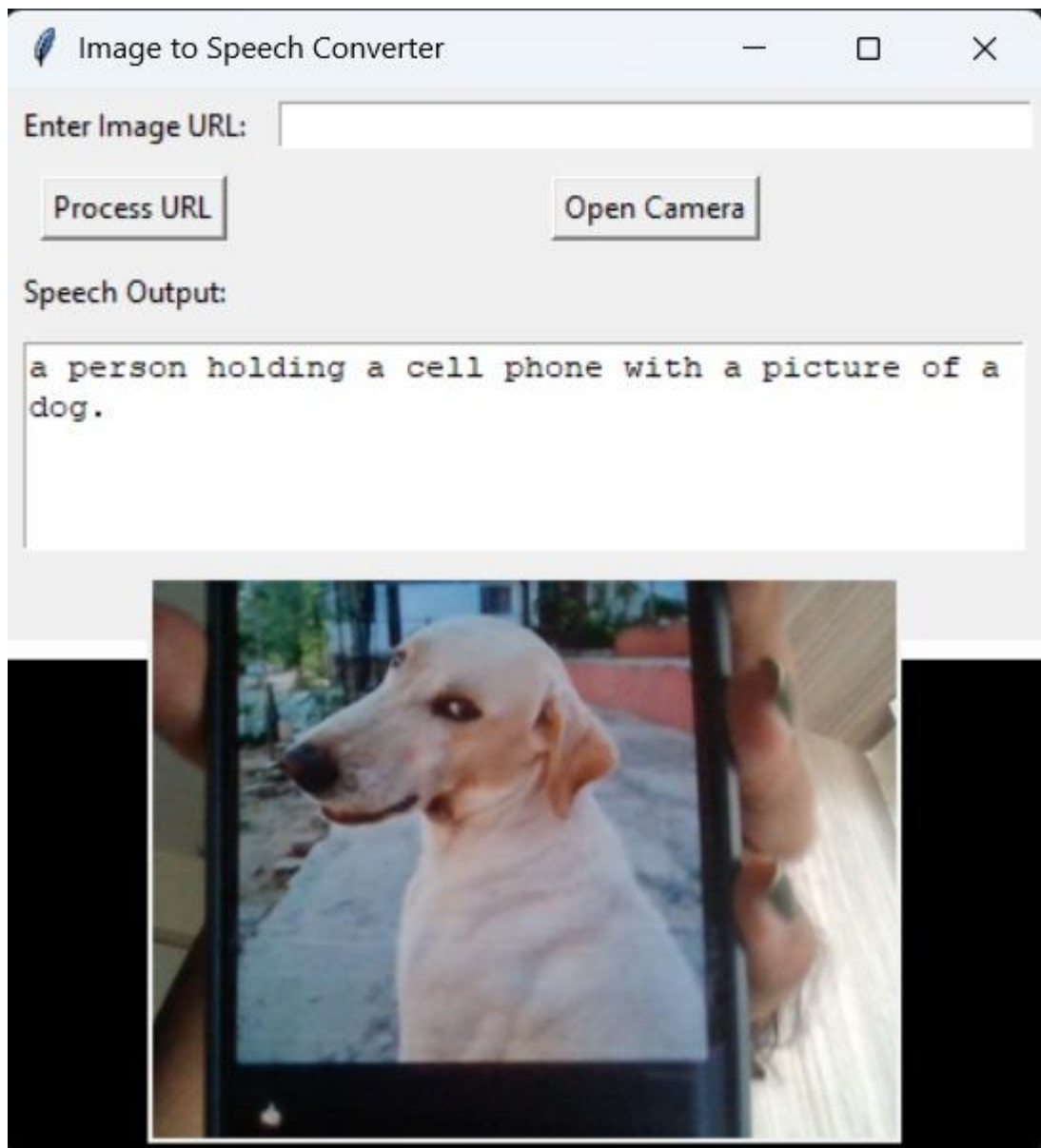


Fig 5.2.1 sample output 1

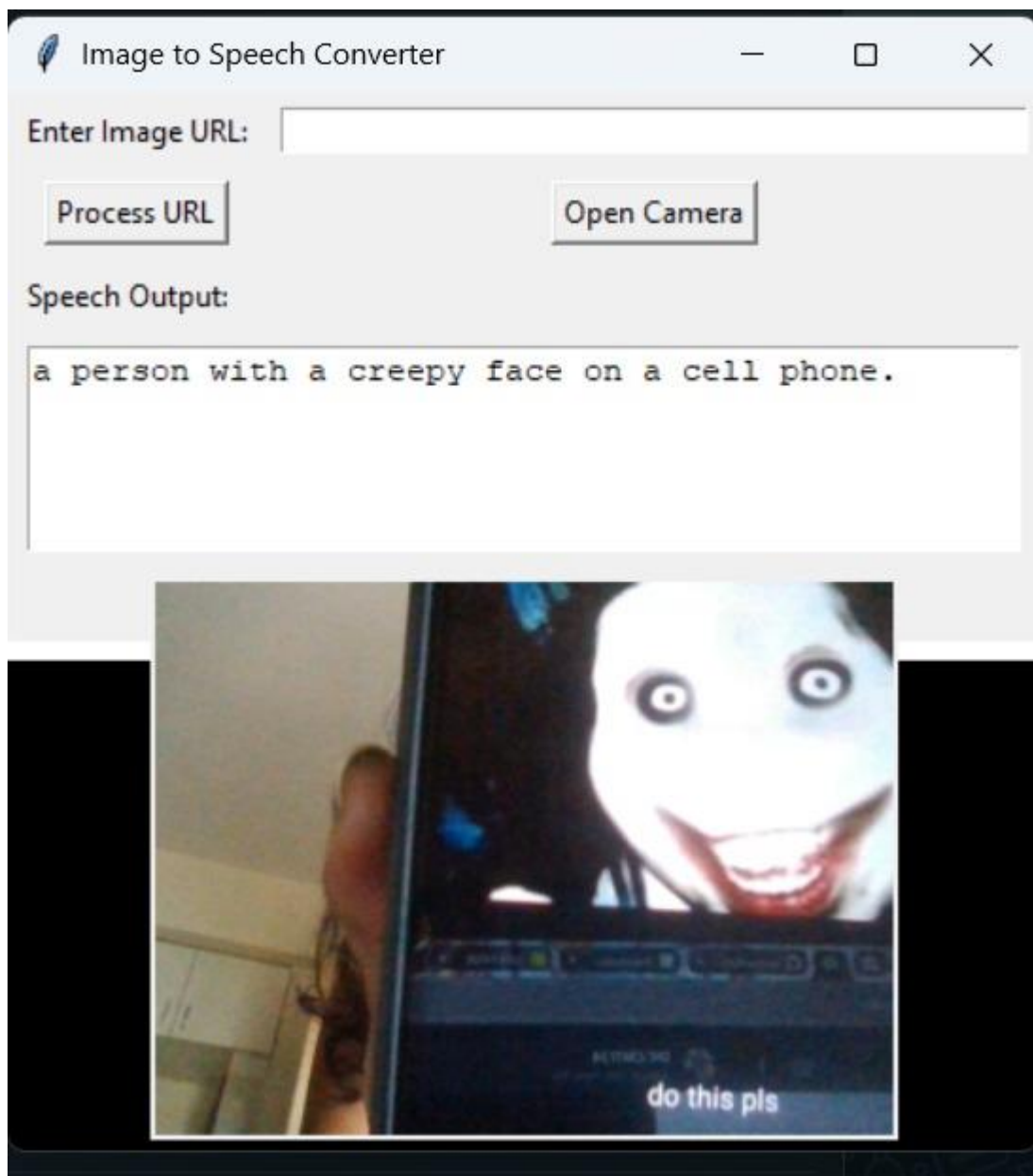


Fig 5.2.2 sample output 2



Real-time camera output transcribed to text enables seamless conversion of visual information into readable text

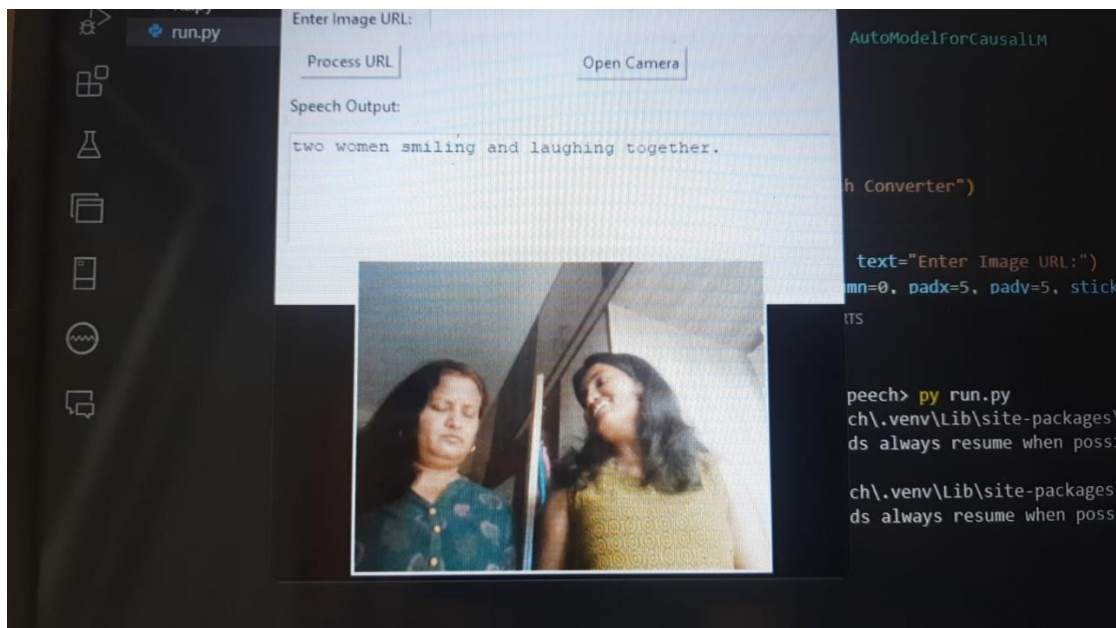


Fig 5.2.3 sample output 3

## CONCLUSION

Our project focuses on the development of a comprehensive system for image-to-speech conversion using cutting-edge machine learning techniques. The ability to convert visual information into spoken language has profound implications for accessibility, education, and assistive technology. Leveraging the rapid advancements in machine learning and computer vision, our system comprises two main components: image-to-text conversion and text-to-speech synthesis. For image-to-text conversion, we employ a Generative Image-to-text Transformer (GIT) model, streamlining the architecture into a single image encoder and text decoder. This model has been trained on extensive data to accurately generate textual descriptions from input images. Additionally, we implement text-to-speech (TTS) synthesis to convert these textual descriptions into spoken language, facilitating seamless accessibility for visually impaired individuals or those preferring auditory information presentation. To further enhance the system's capabilities, we integrate camera access, enabling real-time image-to-speech conversion. Users can capture images directly through their device's camera and receive instant spoken descriptions. Moreover, we introduce a reinforcement learning (RL) model to continually refine the accuracy and relevance of generated text, improving the overall performance of the system.

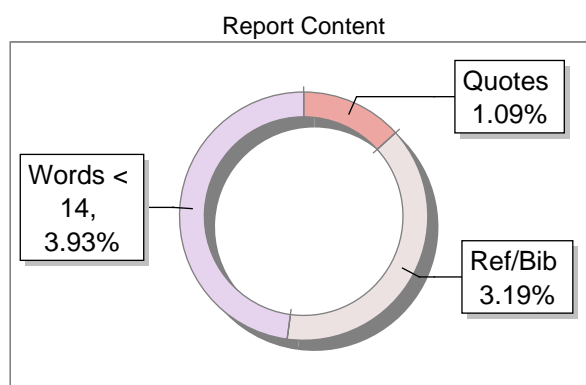
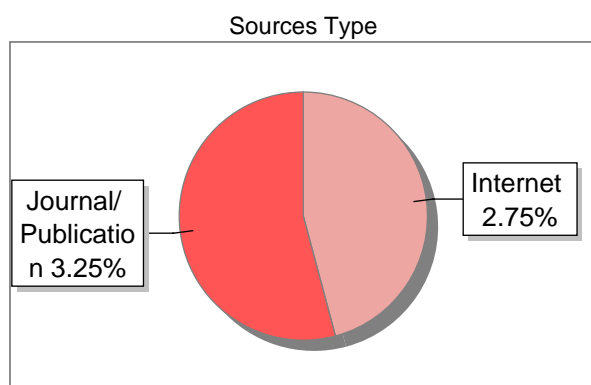
## REFERENCES

- [1] Dr.Sujatha. K, Shruti P. Pati, "Text and Speech Recognition for Visually Impaired People using Optical Character Recognition (OCR), text-to-speech synthesizer (TTS)", IJCRT | Volume 9, (2021).
- [2] Pooja Bendale<sup>1</sup>, Sanika Badhe, Sarthak Bhagat "Text Extraction from Image, Word, PDF and Text-to-Speech Conversion", IRJET| volume 9.(2022).
- [3] Rozelle Jain, Shantanu Das "Image to speech conversion using digital image processing" International Journal of Advanced Research in Engineering and Technology (IJARET), volume 9, (2020)
- [4] Aman Raj Singh, Prof. Diwakar Bhardwaj, Mridal Dixit, Lalit kumar , "An Integrated Model for Text to Text, Image to Text and Audio to Text Linguistic Conversion using Machine Learning Approach" (Li et al., 2023), 2023 6th International Conference on Information Systems and Computer Networks (ISCON) GLA University.
- [5] M.Pandu Babu<sup>1</sup>, Anitha G2 , "OCR Based Image Text to Speech Conversion using K-Nearest Neighbors and Comparing with Fuzzy K-Means Clustering Algorithm", 2023 International Conference on Advances in Computing, Communication and Applied Informatics (ACCAI).
- [6] Mr.Padmavathi P, Bunny Bharadwaj Mahadas ,Shyam Sundhar Kalluri,Praveen Devarapu,Sowndarya Lakshmi Bandi,"Optical Character Recognition and Text to Speech Generation System using Machine Learning ", 2023 2nd International Conference on Applied Artificial Intelligence and Computing (ICAAIC).

### Submission Information

Author Name	MOHITH GOWDA BM
Title	IMAGE TO SPEECH CONVERSION USING MACHINE LEARNING
Paper/Submission ID	1801616
Submitted by	chikmathbasavaraj@gmail.com
Submission Date	2024-05-14 12:07:47
Total Pages, Total Words	39, 6080
Document type	Project Work

### Result Information

Similarity **6 %**

### Exclude Information

Quotes	Excluded
References/Bibliography	Excluded
Source: Excluded < 14 Words	Excluded
Excluded Source	<b>0 %</b>
Excluded Phrases	Not Excluded

### Database Selection

Language	English
Student Papers	Yes
Journals & publishers	Yes
Internet or Web	Yes
Institution Repository	Yes

A Unique QR Code use to View/Download/Share Pdf File





## DrillBit Similarity Report

6

SIMILARITY %

4

MATCHED SOURCES

A

GRADE

A-Satisfactory (0-10%)

B-Upgrade (11-40%)

C-Poor (41-60%)

D-Unacceptable (61-100%)

LOCATION	MATCHED DOMAIN	%	SOURCE TYPE
1	towardsdatascience.com	2	Internet Data
2	cse.anits.edu.in	2	Publication
3	Thesis Submitted to Shodhganga Repository	1	Publication
4	medium.com	1	Internet Data