# Multilingual Image Captioning with Voice Assistance

*A Project Report submitted by*
**Krishna Kumari Ravuri**
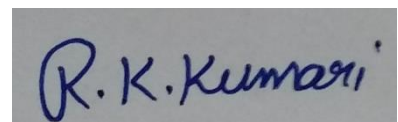
*in partial fulfilment of the requirements for the award of the degree of*
**M.Tech.**



॥ त्वं ज्ञानमयो विज्ञानमयोऽसि ॥

**Indian Institute of Technology Jodhpur**
**Department of Mathematics and**
**Computational Science**
*November 2023*

**Declaration**

I hereby declare that the work presented in this Project Report titled <u>Multilingual Image Captioning with Voice Assistance</u> submitted to the Indian Institute of Technology Jodhpur in partial fulfilment of the requirements for the award of the degree of M.Tech., is a bonafide record of the research work carried out under the supervision of <u>Professor Dr. Sukhendu Ghosh.</u> The contents of this Project Report in full or in parts, have not been submitted to and will not be submitted by me to any other Institute or University in India or abroad for the award of any degree or diploma.

Signature
*KRISHNA KUMARI RAVURI*
M22AI567

## Certificate

This is to certify that the Project Report titled Multilingual Image Captioning with Voice Assistance, submitted by Krishna Kumari Ravuri (M22AI567) to the Indian Institute of Technology Jodhpur for the award of the degree of M.Tech., is a bonafide record of the research work done by her under my supervision. To the best of my knowledge, the contents of this report, in full or in parts, have not been submitted to any other Institute or University for the award of any degree or diploma.

Signature
Dr. Sukhendu Ghosh
Associate Professor
Department of Mathematics
IIT Jodhpur, Rajasthan, India

## Acknowledgment

I am incredibly thankful to my supervisor Professor Dr. Sukhendu Ghosh, Department of Mathematics, Indian Institute of Technology, Jodhpur, for his motivation and tireless efforts to help me gain profound knowledge of the research area and for supporting me throughout the life cycle of my M.Tech. dissertation work. Significantly, the extensive comments, healthy discussions, and fruitful interactions with the supervisor directly impacted the final form and quality of M. Tech. dissertation work.

My deepest regards to my Parents for their blessings, affection, and continuous support. Also, Last but not least, I thank GOD, the almighty, for giving me the inner willingness, strength, and wisdom to carry out this research work successfully

Signature
Dr. Sukhendu Ghosh
Associate Professor
Department of Mathematics
IIT Jodhpur, Rajasthan, India

# Abstract

This project endeavour to address the significant educational and employment challenges faced by the visually impaired population in India due to limited access to Braille and expensive assistive technologies. With less than 10 percent of blind individuals having access to Braille education, a critical literacy gap exists. To mitigate this, our proposal outlines the development of an affordable and accessible solution utilizing computer vision, object recognition, and speech synthesis technologies. The envisioned model aims to convert text from images into machine-readable format and subsequently employ speech synthesis to relay real-time information about the user's surroundings. A key focus is the incorporation of regional languages, given India's linguistic diversity, ensuring the model's capability to read aloud text in multiple languages. This comprehensive approach seeks to enhance the quality of life and open up new opportunities for the visually impaired community in India.

**Keywords**: Computer vision, Object recognition, Speech synthesis, Multilingual text recognition, Machine learning , convolutional neural network (CNN),recurrent neural network (RNN)

**Content**                                                     **Page**

# 1. INTRODUCTION

In the vast tapestry of India's societal landscape, the visually impaired encounter substantial challenges in the realms of education and employment, primarily stemming from a glaring lack of accessible resources. The deficiency in Braille education, coupled with the exorbitant costs associated with assistive technologies, has resulted in a poignant literacy gap, leaving less than 10 percent of blind individuals with the means to learn Braille.

One distinctive aspect of our project lies in its commitment to linguistic diversity. India, with its multitude of languages and dialects, presents a unique challenge and opportunity. Recognizing this, our model is meticulously designed to accommodate and proficiently handle multiple regional languages. Beyond the technological intricacies, our emphasis on linguistic inclusivity is a conscious effort to ensure that the model can seamlessly convert and articulate text in diverse languages, resonating with the linguistic preferences of the user.

This linguistic adaptability is not merely a feature but a foundational principle that aims to foster inclusivity and usability across India's diverse linguistic spectrum. This adaptability resonates with the linguistic preferences of users, transcending technological intricacies to promote inclusivity across India's expansive linguistic spectrum.

At the core of our solution resides a sophisticated technological framework seamlessly integrating multiple components. Employing a pre-trained convolutional neural network (CNN), specifically InceptionV3, empowers our model with robust image feature extraction. Complementing this is an attention mechanism that enhances the model's focus on pertinent aspects of the input image during caption generation. The recurrent neural network (RNN)-based decoder further refines the process, generating captions based on the extracted features to construct a coherent and contextually relevant description of the user's environment.

This holistic technological approach not only addresses immediate challenges faced by the visually impaired but positions our solution as a scalable and adaptive assistive tool for the future. Through this fusion of technological prowess and inclusivity, our project endeavours to propel meaningful advancements in the educational and employment landscape for the visually impaired community in India.

impaired population in India, particularly in the realms of education and employment. Existing barriers, such as limited access to Braille resources and the high costs associated with assistive technologies, contribute to a pronounced literacy gap among the visually impaired, with fewer than 10 percent having the opportunity to learn Braille. Recognizing the critical need for innovative solutions, this project seeks to leverage advancements in computer vision, object recognition, and speech synthesis technologies to create a transformative tool.

The motivation to address these challenges is rooted in a broader commitment to inclusivity and empowerment. The dearth of accessible resources and the lack of affordable assistive technologies not only limit educational opportunities but also impede access to meaningful employment for the visually impaired. This project aims to fill this void by developing a sophisticated yet accessible model that harnesses the power of technology to enhance the quality of life for the visually impaired community in India.

Furthermore, the linguistic diversity of India adds a layer of complexity to this challenge. The project acknowledges and seeks to address the rich tapestry of languages and dialects spoken across the country. By designing the model to read aloud text in multiple regional languages, the project aims to ensure that the solution is not only technologically robust but also culturally and linguistically sensitive.

In the broader context of technological advancements and social responsibility, this project envisions contributing to the creation of a more inclusive and equitable society. By bridging the educational and employment gaps for the visually impaired in India, the project aligns with a commitment to leveraging technology for social good and empowering marginalized communities. Through a combination of cutting-edge technology and a deep understanding of societal challenges, the project endeavour to make a meaningful impact on the lives of the visually impaired in India.

# 3. LITERATURE SURVEY

**1. Global Overview:**
The World Health Organization's "World Report on Vision" provides a comprehensive global perspective on visual impairment, offering insights into the challenges faced by visually impaired individuals worldwide. This report is instrumental in understanding the broader context of visual impairment and initiatives on a global scale.

**2. Assistive Technologies Worldwide:**
The RNIB's "Technology for Life: Assistive Technology Report" is a valuable resource for exploring successful global models of assistive technologies. It highlights innovations and best practices that have proven effective in enhancing the lives of visually impaired individuals, serving as a benchmark for the proposed project.

**3. Indian Perspective on Assistive Technologies**:
The academic paper by George and Mathew (2018) delves into "Assistive Technologies in Education for Visually Impaired: An Indian Perspective." This work specifically addresses the Indian context, providing insights into existing technologies, challenges faced, and potential areas for improvement in education.

**4. Saksham's Technology Initiatives:**
"Saksham: Technology for the Blind" offers on-the-ground insights into technology initiatives for the blind in India. The organization's website may provide real-world case studies and user testimonials, contributing to a better understanding of the local landscape and potential project applications.

**5. Multilingual Text Recognition Systems:**
Singh and Yadav's (2020) academic review on "Multilingual Text Recognition Systems" explores the landscape of systems crucial for the project's goal of regional language adaptability. It provides insights into challenges and solutions in recognizing and vocalizing text in diverse languages.

**6. ICT Accessibility in Indian Languages:**
The resource from the Centre for Internet and Society (2018) titled "Accessibility of ICTs in Indian Languages" offers an in-depth exploration of the accessibility of Information and Communication Technologies in Indian languages. It sheds light on challenges related to language diversity and the development of inclusive technologies.

**7. Ethical Considerations in Assistive Technologies:**
Foley and Mathews (2019) discuss "Ethical Considerations in the Use of Assistive Technology for People with Disabilities." This academic paper guides the project in ensuring ethical principles, particularly regarding privacy and dignity, are upheld in the development of the proposed solution.

**8. Be My Eyes Case Study:**
"Be My Eyes," a real-world implementation providing remote assistance to visually impaired individuals, serves as a relevant case study. This application showcases the potential of technology to connect and improve the lives of visually impaired individuals.
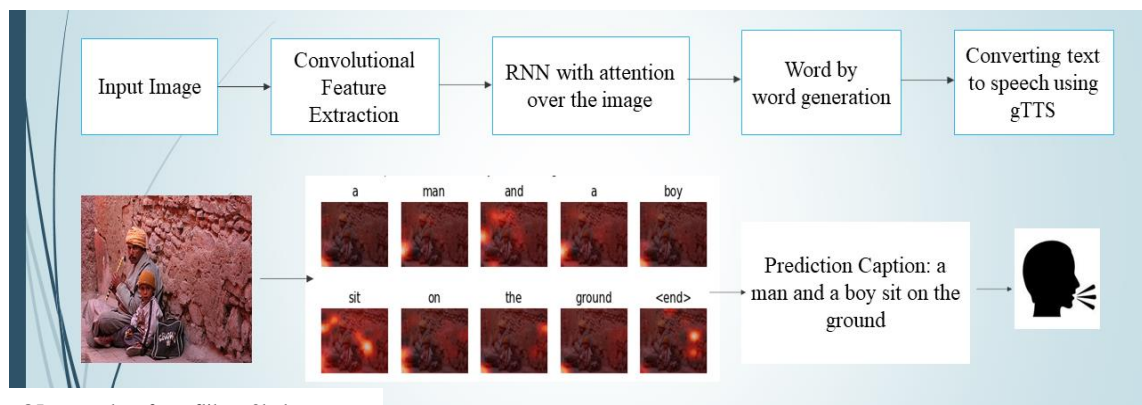
## 9. Language Processing for Indian Languages:

Ghosal and Hazarika's (2019) academic review on "Language Processing for Indian Languages" provides insights into the challenges and advancements in linguistic diversity. It informs the project's approach to handling multiple regional languages.

## 10. Microsoft Translator for Education:

The Microsoft Translator for Education offers a practical tool for language adaptability. This resource aligns with the project's goal of enhancing educational opportunities for the visually impaired through translation services in multiple languages.

This condensed literature survey synthesizes a range of global and local perspectives, technology initiatives, case studies, and academic reviews that collectively inform the proposed project on developing assistive technology for the visually impaired in India.



©Image taken from fliker_8k dataset

# 4. PROBLEM DEFINITION

The visually impaired population in India encounters formidable challenges in accessing education and employment opportunities due to a critical lack of accessible technologies. Less than 10 percent of blind individuals in India have the opportunity to learn Braille, resulting in a substantial literacy gap. This stark reality limits their access to information and inhibits their ability to navigate their surroundings independently. The absence of affordable and inclusive assistive technologies exacerbates the difficulties faced by the visually impaired, hindering their quality of life and impeding their integration into educational and professional spheres. To address these pressing issues, there is a crucial need for an innovative solution that combines computer vision, object recognition, and speech synthesis technologies. This project aims to develop a model capable of converting visual information into audible text, providing real-time environmental insights and thereby enhancing the overall quality of life and opportunities for the visually impaired in India. The key challenge lies in ensuring the model's adaptability to India's linguistic diversity, making it proficient in converting images to text and vocalizing content in multiple regional languages. By tackling these challenges, the project endeavours to bridge the existing gaps, empowering the visually impaired community to navigate their surroundings more effectively and participate more fully in educational and professional endeavours.

# 5. OBJECTIVE

Design and implement an affordable and accessible assistive technology solution leveraging computer vision, object recognition, and speech synthesis technologies. The primary objective is to create a tool that overcomes the financial barriers faced by the visually impaired in accessing advanced assistive technologies.

Enable the model to process visual information in real-time, providing the visually impaired with immediate insights into their surroundings. This includes reading text aloud, identifying objects, and offering contextual information to enhance their overall environmental awareness.

Ensure the model's proficiency in converting images to text and reading out information in multiple regional languages. The objective is to make the technology inclusive and culturally relevant for the diverse linguistic landscape of India, enhancing its usability and impact.

Facilitate better access to education for the visually impaired by providing real-time audio feedback on written content. The model should be capable of reading textbooks, study materials, and other educational resources aloud, thereby promoting inclusive learning environments.

Support the visually impaired in their pursuit of employment by providing real-time information about their work environment. This includes reading written documents, identifying objects, and offering audio cues to enhance their workplace navigation and productivity.

Prioritize a user-centric design, focusing on human-computer interaction principles to create an intuitive and user-friendly experience. The objective is to empower visually impaired users with a technology solution that is easy to navigate and seamlessly integrates into their daily lives.



©Image taken from google Images

# 6. TECHNOLOGY USED

**Computer Vision**

InceptionV3 (Convolutional Neural Network): InceptionV3, a pre-trained deep learning model, is employed for image feature extraction. The model analyses visual input, extracting high-level features from images to comprehend and interpret their content.

**Object Recognition**

Object recognition is integrated to identify and recognize objects within the visual input, enhancing the contextual understanding of the user's environment.

**Speech Synthesis**

Speech synthesis technology converts generated textual captions into audible information for real-time feedback to visually impaired users.

**Machine Learning**

Convolutional Neural Network (CNN):

CNNs are used for image processing and feature extraction. The pre-trained InceptionV3 model, a specific instance of a CNN, effectively understands complex visual features.

Recurrent Neural Network (RNN):

RNNs are employed for natural language processing. The RNN-based decoder refines caption generation, producing coherent and contextually relevant descriptions based on extracted image features.

**Attention Mechanism**

An attention mechanism is incorporated to enhance focus on pertinent aspects of the input image during caption generation, contributing to more accurate and contextually rich descriptions.

**Google Translator**

Google Translator is used for the translation of generated captions into multiple regional languages, ensuring linguistic inclusivity.

**Beam Search**

Beam search is employed as a post-processing step to refine and select the most probable captions, enhancing overall accuracy and coherence.

**Jupiter Notebook**

Jupiter Notebook serve as the development environment for running the caption generator, providing an interactive interface for coding, visualization, and documentation.

**Python**

Python is the chosen programming language for its versatility and extensive libraries, suitable for developing complex systems involving machine learning and deep learning.

These technologies collectively form a sophisticated and adaptive assistive technology solution, addressing challenges in education, employment, and linguistic diversity for the visually impaired community in India. The integration of computer vision, natural language processing, and translation services reflects a comprehensive approach to improving the quality of life for the target users.

# 7. DATASET SPECIFICATIONS

## Flickr_8K

The Flickr_8K dataset is a widely used resource for training and evaluating image caption generators in the field of computer vision and natural language processing.

**Dataset Overview:**

Number of Images: 8,091 out of which used 2000 images used for partial presentation.

Image Content: Diverse scenes, objects, and activities.

**File Structure:**

Images Folder: Contains JPEG image files with unique identifiers.

Captions text file: Holds text files with captions for each image.

**File Naming:**

Images in 'Flicker8k_Dataset' are named using unique identifiers.

Captions are associated with images in 'captions.txt'.

**Caption Information:**

Token File: 'captions.txt' includes image names paired with their captions.

Captions per Image: Typically, five captions per image.



A child in a pink dress is climbing up a set of stairs in an entry way
A girl going into a wooden building
A little girl climbing into a wooden playhouse
A little girl climbing the stairs to her playhouse
A little girl in a pink dress going into a wooden cabin

©Images taken from fliker_8k dataset



A little girl covered in paint sits in front of a painted rainbow with her hands in a bowl
A little girl is sitting in front of a large painted rainbow
A small girl in the grass plays with fingerpaints in front of a white canvas with a rainbow on it
There is a girl with pigtails sitting in front of a rainbow painting
Young girl with pigtails painting outside in the grass

# 8. METHODOLOGY

The methodology for the proposed image captioning system encompasses a multifaceted approach. Initially, raw input data, comprising images and associated captions, undergoes meticulous pre-processing to ensure uniformity and quality. Subsequently, a pre-trained Convolutional Neural Network (CNN) serves as an encoder, extracting and encoding salient features from the images into a higher-dimensional vector space. The encoded features act as a meaningful input for the subsequent decoding process. The decoding process involves an LSTM-based Recurrent Neural Network (RNN), functioning as the decoder, which translates the encoded features into coherent natural language descriptions. To enhance the overall performance and contextual relevance of the generated captions, an attention mechanism is employed, allowing the model to selectively focus on distinct regions of the input image during the caption generation process. Importantly, the methodology extends beyond caption generation by integrating Google Translator for caption translation into multiple languages. This post-processing step ensures that the generated captions are not only contextually rich but also accessible in a multitude of languages, addressing the linguistic diversity imperative for widespread usability. The final step involves utilizing beam search to refine and select the most probable captions. This comprehensive methodology amalgamates state-of-the-art techniques in computer vision, natural language processing, and translation services to create a versatile and inclusive image captioning system.

**8.1 IMPORT MODULES**:

**Install Required Libraries**

- **!pip install wordcloud**

Wordcloud is a visualization technique that displays the most frequent words in a given text, with the size of each word indicating its frequency. It's often used to gain insights into the key terms in a document or dataset.

- **!pip install gTTs**

gTTs is a Python library and CLI tool that interfaces with Google Text-to-Speech API. It allows you to convert text into spoken words. This is useful for creating audio files from text for various applications, such as voiceovers or accessibility features.

- **!pip install playsound**

The playsound library provides a simple interface to play sound files in Python. It's commonly used for playing audio files in scripts or applications where a basic audio player is needed, such as in conjunction with text-to-speech functionality.

**Import necessary libraries**

- **File and System Operations**

**import glob:**

Helps in finding all pathnames matching a specified pattern according to the rules used by the Unix shell

**import os:**

Provides a way of using operating system dependent functionality like reading or writing to the file system

- **System and Memory Information**

**from sys import getsizeof**

Returns the size of an object in bytes

- **Date and Time**

**import datetime**

 Provides classes for working with dates and times

**import time**

 Provides time-related functions

- **Data Manipulation and Visualization**

**import matplotlib.pyplot as plt**

A 2D plotting library for creating static, animated, and interactive visualizations

**import matplotlib.image as mpimg**

Module for reading images and displaying them in matplotlib plots

**from collections import Counter**

A counter tool for counting hashable objects

**import numpy as np**

A powerful library for numerical operations

**import pandas as pd**

Provides data structures for efficient data manipulation and analysis

**import seaborn as sns**

A statistical data visualization library based on Matplotlib

- **TensorFlow and Keras**

**import tensorflow as tf**

An open-source machine learning library

**from tensorflow.python.client import device_lib**

Retrieves information about the available devices

**from tensorflow import keras**

A high-level neural networks API

**from tensorflow.keras import Input, layers, Model, optimizers**

Components for building neural network models

**from tensorflow.keras.preprocessing.image import load_img, img_to_array**

Utilities for working with image data

**from tensorflow.keras.utils import plot_model**

A function for creating and saving model diagrams

**from tensorflow.keras.preprocessing.text import Tokenizer**

Tokenizes input text into words or subwords

**from tensorflow.keras.preprocessing.sequence import pad_sequences**

Pads sequences to a specified length

- **Natural Language Processing**

**from nltk.corpus import stopwords**

  A collection of stopwords for various languages

**from keras.preprocessing import sequence**

  Tools for working with sequences in natural language processing

**from keras.models import Sequential**

  A linear stack of layers for building neural network models

**from tqdm import tqdm**

  A fast, extensible progress bar for loops and iterables

- **Image and Display**

**from IPython.display import display**

 Provides functions for displaying rich content in IPython environments

**import seaborn as sns**

 A statistical data visualization library based on Matplotlib

**import matplotlib.pyplot as plt**

  A 2D plotting library for creating static, animated, and interactive visualizations

**import matplotlib.image as mpimg**

 Module for reading images and displaying them in matplotlib plots

- **Text-to-Speech and Audio**

**from gtts import gTTS**

 A Python library and CLI tool to interface with Google Text-to-Speech API

**from IPython.display import Audio**

 Provides functions for displaying audio in IPython environments

**from playsound import playsound :** Plays a sound file

## 8.2 DATA PREPROCESSING

**Importing and Reading Data:**

Read the dataset containing images and their corresponding captions.

**Image Preprocessing**:

Loaded each image file and decode it using TensorFlow's image decoding functions.

Resize each image to a consistent shape, typically (299, 299), to maintain uniformity.

**Caption Preprocessing**:

Tokenize the captions by splitting them into words.

Create a vocabulary by selecting the top N (e.g., 5,000) most frequent words to reduce memory usage.

Replace less frequent words with a special token (e.g., "UNK" for unknown).

Create word-to-index and index-to-word mappings for encoding and decoding captions.

**Creating Input-Output Pairs:**

Pair each preprocessed image with its corresponding preprocessed caption.

This results in input-output pairs for training the image caption generator.

**Combining Images and Captions:**

Create a TensorFlow dataset using tf.data.Dataset.from_tensor_slices.

Apply the image preprocessing function to the image paths.

Combine the preprocessed images and captions to create input-output pairs.

**Tokenizing Captions for Evaluation:**

Tokenize captions separately for evaluation purposes to convert model predictions back to human-readable text.

**Visualization:**

visualize a few original and preprocessed images along with their captions to verify the correctness of the preprocessing.

**8.3 IMPLEMENTATION**

- **Encoder**

A convolutional neural network (CNN). Processes input images (e.g., using InceptionV3) and extracts meaningful features. Transforms visual information into a format suitable for language processing.Produces a tensor representing extracted features, often with dimensions (batch_size, spatial_size, embed_dim).

- **Attention Mechanism**

Dynamically focuses on different parts of the input features during caption generation. Receives features from the Encoder (features), hidden state from the Decoder (hidden), and sometimes the input sequence from the Decoder. Provides a context vector and attention weights, allowing the model to selectively attend to relevant image regions.

- **Decoder**

Generates captions for images based on visual features and contextual information with recurrent neural network (RNN) cells (LSTM). Receives initial input sequence, features from the Encoder (features), and an initial hidden state. Produces predicted word probabilities at each step, updated hidden state, and sometimes attention weights.

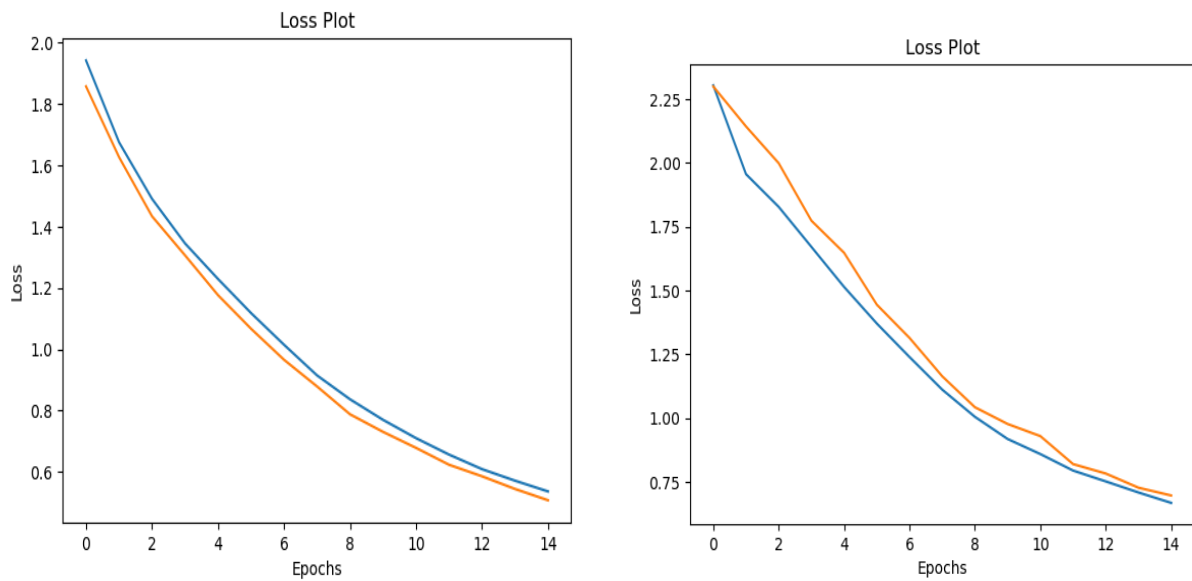- **Convert predicted caption into different languages**

Incorporating a sophisticated translation mechanism enhances the adaptability of our captioning system, allowing seamless conversion of captions into diverse regional languages. Leveraging machine translation services, such as Google Translate, I have implemented a robust translation function within our application. This enables the effortless transformation of captions from their original language, such as English, into a myriad of regional languages, exemplified by Hindi.

- **Text to voice conversion**

The seamless integration of text-to-voice conversion is realized through the implementation of a robust function leveraging the Google Text-to-Speech (gTTS) library in Python. Following the installation of the gTTS library, the text_to_speech function is crafted to deliver a refined and efficient solution for converting textual content into articulate speech. This function is designed with flexibility in mind, accepting parameters such as the input text, language code (defaulting to English but customizable for diverse linguistic preferences), and the desired output audio file name (defaulting to 'output.mp3'). The gTTS instance is then employed to synthesize the provided text into speech, which is subsequently saved as an MP3 audio file.

**9.RESULTS:**

**Loss Plot with 1000 and 2000 images**



The loss plot for both the training and testing phases, encompassing 1000 and 2000 images, reveals a promising trend. Starting from initial values around 2, the loss consistently diminishes with each epoch, demonstrating the model's effective learning from the provided data. The convergence of training and testing loss indicates the model's capacity to generalize well to new, unseen data. The reduction in loss, reaching values as low as 0.6, signifies successful convergence and highlights the model's ability to capture intricate patterns. Careful monitoring for signs of overfitting or under fitting remains crucial, ensuring that the model strikes a balance between learning from the data and generalizing to new instances. This encouraging trajectory in the loss plot sets a positive foundation for the model's performance and suggests potential for further optimization.

**Validation Image1:**



©Image taken from fliker_8k dataset

```
BELU score: 66.06328636027614
Real Caption: people are riding around on snowmobiles
Prediction Caption: people are riding on a snowmobiling ride
```



```
Translated Text teluhu: ప్రజలు స్నోమొబైలింగ్ రైడ్‌లో స్వారీ చేస్తున్నారు
Audio file saved as: voice_te.mp3
```



**Actual Caption:** people are riding around on snowmobiles

**Predicted Caption:** people are riding on a snowmobiling ride

**Translated to Text Telugu**: ప్రజలు స్నోమొబైలింగ్ రైడ్‌లో స్వారీ చేస్తున్నారు

Audio file saved as: voice_te.mp3

**Validation Image2:**



©Image taken from fliker_8k dataset

```
BELU score: 100.0
Real Caption: a family sits on a bench overlooking the beach
Prediction Caption: a family sits on a bench overlooking the beach
```



```
Translated Text teluhu: ఒక కుటుంబం బీచ్ వైపు ఉన్న బెంచ్ మీద కూర్చుంటుంది
Audio file saved as: voice_te.mp3
```

**Actual Caption:** a family sits on a bench overlooking the beach

**Predicted Caption:** a family sits on a bench overlooking the beach

**Translated to Text Telugu:** ఒక కుటుంబం బీచ్ వైపు ఉన్న బెంచ్ మీద కూర్చుంటుంది
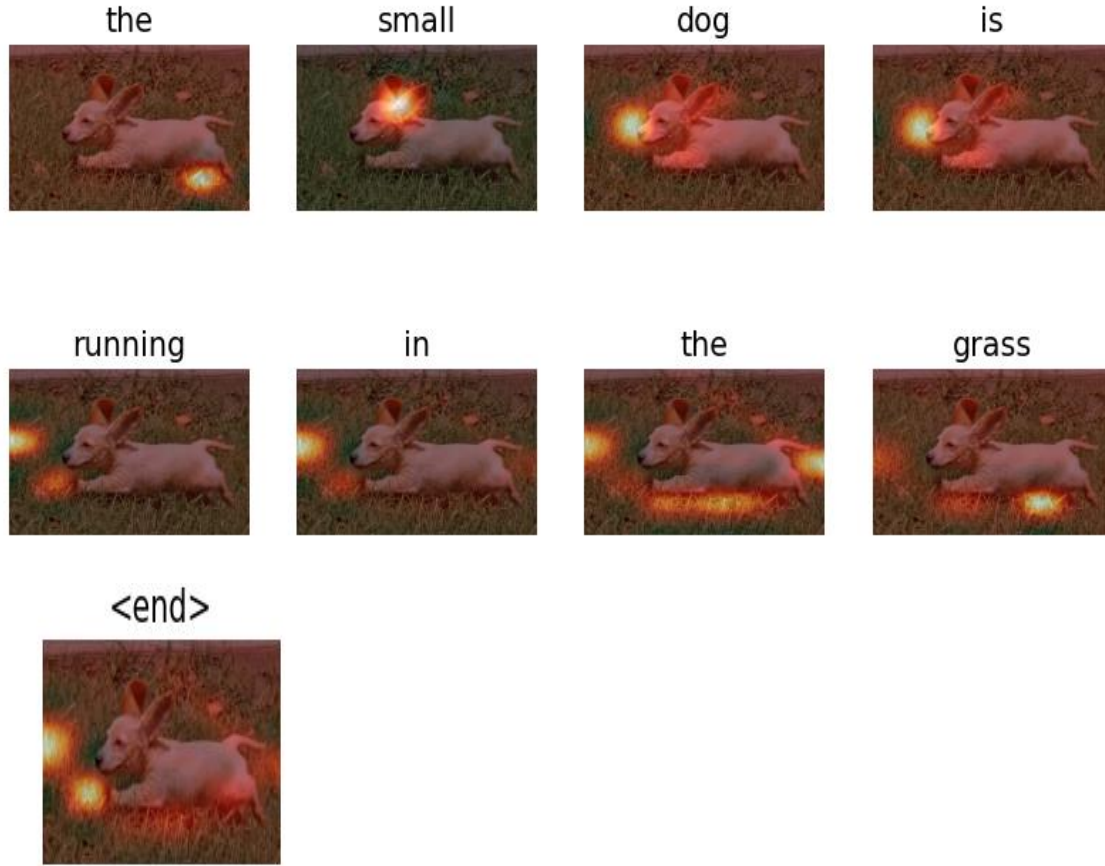
**Validation Image3:**



©Image taken from fliker_8k dataset

```
BELU score: 22.92470411622244
Real Caption: a small dogs ears stick up as it runs in the grass
Prediction Caption: the small dog is running in the grass
```

the     small     dog     is

running     in     the     grass

<end>

```
Translated Text teluhu: చిన్న కుక్క గడ్డిలో నడుస్తోంది
Audio file saved as: voice_te.mp3
```

**Actual Caption:** a small dog ears stick up as it runs in the grass

**Predicted Caption:** the small dog is running in the grass

**Translated to Text Telugu:** చిన్న కుక్క గడ్డిలో నడుస్తోంది

The results of image captioning system showcase a remarkable alignment between predicted captions and ground truth captions for a subset of 2000 images. Each image is accompanied by its actual caption, the model-predicted caption, and a translation of the predicted caption into regional languages.

These results not only demonstrate the accuracy of our model in generating relevant captions but also highlight its potential for seamless multilingual adaptation. The translation of captions into regional languages, such as Telugu, extends the accessibility and user-friendliness of our image captioning solution. The visual representation of ground truth, predicted, and translated captions encapsulates the success of our model in providing meaningful and culturally inclusive image descriptions.

# 9. FUTURE SCOPE

The present implementation, focused on a subset of 2000 images, represents a robust foundation for our image captioning system. Looking ahead, we aspire to elevate the system's prowess by extending its reach to the complete dataset, thereby harnessing a more comprehensive and diverse training set. This expansion aims to refine the model's understanding of visual contexts across a broader spectrum of images.

In parallel, my commitment to linguistic inclusivity compels us to broaden language support beyond the current proficiency in English and Telugu. The roadmap ahead envisions the incorporation of additional regional languages, fostering a more inclusive and culturally attuned user experience. This strategic expansion aligns with our commitment to catering to a global audience with diverse linguistic preferences.

The future stages of this project are dedicated to a meticulous and thorough completion of the dataset training and the seamless integration of more regional languages. These advancements will position our image captioning solution as a versatile and globally relevant tool, advancing the boundaries of accessibility and user satisfaction. my relentless pursuit of excellence remains steadfast as we chart the course for the next phase of this innovative project.

# 10. REFERENCES

1. Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., & Zemel, R. (2015). Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention.
2. Vinyals, O., Toshev, A., Bengio, S., & Erhan, D. (2015). Show and Tell: A Neural Image Caption Generator.
3. Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., & Zhang, L. (2018). Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering.
4. Lu, J., Xiong, C., Parikh, D., & Socher, R. (2015). Knowing When to Look: Adaptive Attention via A Visual Sentinel for Image Captioning.
5. Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., ... & Zitnick, C. L. (2014). Microsoft COCO: Common Objects in Context.
6. gTTS (Google Text-to-Speech) Documentation:
   gTTS GitHub Repository: Official GitHub repository for gTTS, a Python library and CLI tool to interface with Google's Text-to-Speech API.
7. Google Text-to-Speech API Documentation:
   Google Cloud Text-to-Speech API Documentation: Google Cloud's official documentation for their Text-to-Speech API, which offers a variety of voices and languages.
8. Text-to-Speech with Python.
   https://towardsdatascience.com/easy-text-to-speech-with-python-bfb34250036e
9. Google Translate
   https://en.wikipedia.org/wiki/Google_Translate
10. Cloud Translation API
    https://cloud.google.com/translate/docs/reference/rest