# VisionEcho Project:

## Image to Speech for the Blind



Group Members:

Tayyab Bin Qamar (CR-22041)

Maryam Khan (CR-22021)

Ayesha Yousuf (CR-22004)

## Exclusive Summary

*The VisionEcho project is a smart system that helps blind and visually impaired people understand what is shown in an image. It does this in two main steps: first, it looks at a picture and creates a sentence describing what's happening (like "a man is riding a bike on the street"). Then, it reads that sentence out loud using a voice tool. This means someone who can't see the image can still know what's in it just by listening.*

*We use a special deep learning model to do this. It includes a CNN (InceptionV3) to see and understand the image, and an RNN (GRU) to write a caption about the image. We also use an attention mechanism so the model knows which part of the image to focus on when writing each word. Finally, we use a Text-to-Speech (TTS) tool to speak the caption.*

*This project was built using Python with TensorFlow and other helpful libraries like NumPy, Pandas, and gTTS. The data used for training came from the Flickr8k dataset, which includes thousands of images and their descriptions. The final goal of this project is to make visual information more accessible for people who cannot see.*

## 1. Introduction

*The **VisionEcho** project aims to assist visually impaired individuals by converting visual content into descriptive audio. According to the World Health Organization (WHO), approximately **285 million people** suffer from visual impairments, with **39 million** being completely blind. This project leverages **Deep Learning (CNN-RNN with Attention Mechanism)** and **Natural Language Processing (NLP)** to generate textual descriptions of images, which are then converted into speech using a **Text-to-Speech (TTS)** API.*

## 2. Problem Statement

*Visually impaired individuals face difficulties navigating visual content, especially on social platforms. Although some platforms offer basic alt-text, they lack detailed scene understanding. VisionEcho bridges this gap by generating descriptive image captions and translating them into speech, offering a more natural and rich interaction with visual media.*

## *3. Methodology*

*The project follows a **two-step approach**:*

1. ***Image-to-Caption Generation*** *(CNN-RNN with Attention Mechanism)*
2. ***Caption-to-Audio Conversion*** *(Text-to-Speech API)*

### *3.1. Dataset Used*

*The **Flickr8k dataset** was used, containing **8,091 images** and **40,455 captions**. Each image has multiple human-annotated captions for training.*

### *3.2. Architecture Overview*

*The system consists of:*

- ***Encoder (CNN-based)**: Extracts image features.*
- ***Decoder (RNN-based)**: Generates captions using attention.*
- ***Text-to-Speech (TTS)**: Converts captions into audio.*

## *3.3. Overview of Libraries and Algorithms*

*Libraries Used*

1. ***TensorFlow (tensorflow)** – Core deep learning framework used for model creation, training, and inference.*
2. ***NumPy (numpy)** – Supports mathematical computations and matrix operations.*
3. ***Pandas (pandas)** – Handles data manipulation for caption preprocessing.*
4. ***Matplotlib (matplotlib.pyplot)** – Visualizes images and attention maps.*
5. ***Scikit-learn (sklearn.model_selection)** – Splits data into training and testing sets.*
6. ***glob** – Retrieves file paths matching patterns (e.g., image files).*
7. ***os** – Ensures platform-independent file and directory path handling.*
8. ***gTTS (gtts)** – Converts generated text captions into speech audio.*
9. ***IPython.display** – Enables audio playback inside Jupyter Notebooks.*

*Algorithms and Models Used*

1. **InceptionV3 (CNN)**
   - *A pre-trained convolutional neural network from TensorFlow's Keras library.*
   - *Used for image feature extraction (encoder part).*
   - *include_top=False removes the classification layer, yielding feature maps.*
2. **Attention Mechanism**
   - *Enables the model to dynamically focus on different parts of the image during caption generation.*
   - *Provides a context vector that guides word prediction based on visual relevance.*
3. **GRU (Gated Recurrent Unit)**
   - *A variant of RNN used in the decoder.*
   - *Retains memory over time steps, handling sequential dependencies in language.*
   - *More efficient than LSTM while maintaining performance.*
4. **Tokenizer and Sequence Padding**
   - *Converts text to integer sequences (tokenization).*
   - *Pads sequences to ensure uniform input size for batch training.*
5. **Text-to-Speech (gTTS)**
   - *Transforms generated captions into audible speech.*
   - *Provides a human-like voice output for real-time user interaction*

# 4. Implementation Steps

## 4.1. Importing Required Libraries

```
import glob

import matplotlib.pyplot as plt

import numpy as np

import pandas as pd

import tensorflow as tf

from sklearn.model_selection import train_test_split

import os
```

*Explanation:*

- *glob: Retrieves pathnames matching a pattern.*
- *matplotlib.pyplot: Used to display and visualize images.*
- *numpy: Supports numerical operations and matrix computations.*
- *pandas: Handles dataset manipulation and analysis.*
- *tensorflow: Framework used to build and train the deep learning model.*
- *sklearn.model_selection: Splits the dataset into training and test sets.*
- *os: Enables directory and file path manipulations.*

*Without these foundational libraries, the system would lack key functionalities needed to manage data, process text and images, and construct and train models. Each library plays a distinct role, forming the backbone of the project.*

## 4.2. Loading and Preprocessing the Dataset

```
DATA_DIR = '/home/dl_content/Flickr Dataset'

IMAGES_DIR = os.path.join(DATA_DIR, 'Images')

CAPTIONS_FILE = os.path.join(DATA_DIR, 'captions.txt')



all_imgs = glob.glob(IMAGES_DIR + '/*.jpg')
print("Total images:", len(all_imgs))
```

*Explanation:*

- *Defines paths for images and captions.*
- *glob loads all image files with .jpg extension.*
- *os.path.join ensures compatibility across operating systems..*

*This stage is vital to ensure the system has access to all image files it needs for training. If improperly loaded, the model would not have the necessary visual input to learn caption generation.*

## 4.3. Text Tokenization and Padding

```
tokenizer = tf.keras.preprocessing.text.Tokenizer(num_words=5000, oov_token="<unk>")

tokenizer.fit_on_texts(annotations)
```

```
train_seqs = tokenizer.texts_to_sequences(annotations)

max_sequence_len = max(len(t) for t in train_seqs)
cap_vector = tf.keras.preprocessing.sequence.pad_sequences(train_seqs, padding='post',
maxlen=max_sequence_len)
```

**Explanation**:

- **Tokenizer:** Converts words to unique integer tokens.
- **fit_on_texts:** Learns word frequencies.
- **texts_to_sequences:** Converts each caption into a list of tokens.
- **pad_sequences:** Ensures equal-length sequences for input to the neural network.

*Neural networks require fixed-size inputs. Tokenization transforms natural language into numbers the model can learn from, while padding ensures consistent input size across all samples.*

## 4.4. Image Feature Extraction using InceptionV3

```
image_model = tf.keras.applications.InceptionV3(include_top=False, weights='imagenet')

new_input = image_model.input

hidden_layer = image_model.layers[-1].output
image_features_extract_model = tf.keras.Model(new_input, hidden_layer)
```

**Explanation:**

- **InceptionV3:** A pre-trained CNN used to extract image features.
- **include_top=False:** Removes the final classification layer.
- **image_features_extract_model:** Outputs the feature map instead of classification.
- Utilizes a pre-trained InceptionV3 CNN to extract high-level visual features.
- Removes the final classification layer

*This CNN processes images and outputs a feature map, which is later used to generate context-aware captions. Using a pre-trained model accelerates development and improves accuracy.*

## 4.5. Encoder-Decoder Model with Attention

### 4.5.1. Encoder

```python
class Encoder(tf.keras.Model):

    def __init__(self, embed_dim=256):

        super(Encoder, self).__init__()

        self.dense = tf.keras.layers.Dense(embed_dim)


    def call(self, features):

        features = self.dense(features)

        features = tf.nn.relu(features)
        return features
```

**Explanation**:

- Projects high-dimensional CNN features to a lower-dimensional embedding.
- Applies ReLU for non-linear transformation.
- Transforms raw image features into a dense embedding vector.

This compact representation makes it easier for the decoder to interpret complex image data. ReLU introduces non-linearity, helping the model learn intricate patterns.

### 4.5.2. Attention Mechanism

```python
class Attention_model(tf.keras.Model):

    def __init__(self, units):

        super(Attention_model, self).__init__()

        self.W1 = tf.keras.layers.Dense(units)

        self.W2 = tf.keras.layers.Dense(units)
```

```
    self.V = tf.keras.layers.Dense(1)


  def call(self, features, hidden):

    hidden_with_time_axis = tf.expand_dims(hidden, 1)

    attention_hidden_layer = tf.nn.tanh(self.W1(features) + self.W2(hidden_with_time_axis))

    score = self.V(attention_hidden_layer)

    attention_weights = tf.nn.softmax(score, axis=1)

    context_vector = attention_weights * features

    context_vector = tf.reduce_sum(context_vector, axis=1)
    return context_vector, attention_weights
```

**Explanation**:

- Computes **attention weights** to focus on relevant image regions.
- Generates a **context vector** for the decoder.
- Calculates attention weights to focus on specific parts of the image while generating each word.

Attention ensures that the decoder pays more focus to relevant areas of the image depending on the word it is generating. This mimics human visual processing and significantly improves caption quality.

## 4.5.3. Decoder (RNN with GRU)

```
class Decoder(tf.keras.Model):

  def __init__(self, embed_dim, units, vocab_size):

    super(Decoder, self).__init__()

    self.attention = Attention_model(units)

    self.embed = tf.keras.layers.Embedding(vocab_size, embed_dim)

    self.gru = tf.keras.layers.GRU(units, return_sequences=True, return_state=True)
```

```python
    self.d1 = tf.keras.layers.Dense(units)

    self.d2 = tf.keras.layers.Dense(vocab_size)


  def call(self, x, features, hidden):

    context_vector, attention_weights = self.attention(features, hidden)

    x = self.embed(x)

    x = tf.concat([tf.expand_dims(context_vector, 1), x], axis=-1)

    output, state = self.gru(x)

    x = self.d1(output)

    x = tf.reshape(x, (-1, x.shape[2]))

    x = self.d2(x)
    return x, state, attention_weights
```

**Explanation**:

- Uses **GRU (Gated Recurrent Unit)** for sequential caption generation.
- Combines **context vector** and **word embeddings** for prediction.

- Uses word embeddings, GRU cells, and dense layers to sequentially generate words.

This decoder takes in the image context and the previous word to predict the next word. GRU is chosen for its efficiency and ability to retain memory over time, essential for language modeling.

## 4.6. Training the Model

```python
optimizer = tf.keras.optimizers.Adam()

loss_object = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)


def loss_function(real, pred):

  mask = tf.math.logical_not(tf.math.equal(real, 0))
```

```python
        loss_ = loss_object(real, pred)

    mask = tf.cast(mask, dtype=loss_.dtype)

    loss_ *= mask

    return tf.reduce_mean(loss_)


@tf.function
def train_step(img_tensor, target):

    loss = 0

    hidden = decoder.init_state(batch_size=target.shape[0])

    dec_input = tf.expand_dims([tokenizer.word_index['<start>']] * target.shape[0], 1)

    with tf.GradientTape() as tape:

        features = encoder(img_tensor)

        for i in range(1, target.shape[1]):

            predictions, hidden, _ = decoder(dec_input, features, hidden)

            loss += loss_function(target[:, i], predictions)

            dec_input = tf.expand_dims(target[:, i], 1)

    gradients = tape.gradient(loss, trainable_variables)

    optimizer.apply_gradients(zip(gradients, trainable_variables))
    return loss
```

**Explanation**:

- **Adam optimizer** minimizes the loss.
- **Teacher Forcing** improves training stability.

- Handles the training of the encoder-decoder model using teacher forcing and gradient descent.

*This function ensures efficient backpropagation and updates model weights to minimize captioning errors. The masking step avoids computing loss for padded tokens*

## 4.7. Evaluation (BLEU Score & Attention Visualization)

```python
def evaluate(image):

    hidden = decoder.init_state(batch_size=1)

    img_tensor_val = image_features_extract_model(tf.expand_dims(preprocess_image(image)[0], 0))

    features = encoder(img_tensor_val)

    dec_input = tf.expand_dims([tokenizer.word_index['<start>']], 0)

    result = []

    for i in range(max_sequence_len):

        predictions, hidden, _ = decoder(dec_input, features, hidden)

        predicted_id = tf.argmax(predictions, 1).numpy()[0]

        result.append(tokenizer.index_word[predicted_id])

        if tokenizer.index_word[predicted_id] == '<end>':

            break

        dec_input = tf.expand_dims([predicted_id], 0)
    return ' '.join(result)
```

**Explanation**:

- **Greedy Search** generates captions word-by-word.
- **BLEU Score** evaluates caption quality.

- Evaluates a single image and generates a caption using greedy decoding.

*This inference step is critical for testing the model on new, unseen images. It reflects the model's ability to generalize and narrate meaningful descriptions.*

### 4.8. Text-to-Speech Conversion

```python
from gtts import gTTS

import IPython.display as display


soundFile = 'pred_caption.mp3'

tts = gTTS("A dog is playing in the park", slow=False)

tts.save(soundFile)
display.display(display.Audio(soundFile))
```

**Explanation**:

- **Google Text-to-Speech (gTTS)** converts text into speech.

- Converts the generated caption into audio and plays it in the notebook.

Text-to-speech technology bridges the final accessibility gap, converting machine-generated captions into audible speech so visually impaired users can benefit in real time.

## 5. Results

- **BLEU Score**: Achieved **0.65** (1-gram), indicating good caption accuracy.
- **Attention Visualization**: Model correctly focuses on relevant objects.
- **Real-time Audio Generation**: Successfully converts captions into speech.

The VisionEcho model achieved reasonable accuracy in describing images and consistently emphasized relevant visual features through attention maps. The auditory feedback component made the system usable for its intended audience, demonstrating real-world feasibility.

## *6. Conclusion*

*VisionEcho is a promising tool that translates the visual world into descriptive audio for the blind. Its layered architecture using CNN-RNN with attention and TTS synthesis has shown reliable performance in captioning and narration. Future work should focus on using larger datasets (e.g., MS-COCO), integrating transformer-based models, and deploying lightweight versions on mobile devices for real-time accessibility support.*

*This project proves that deep learning can meaningfully augment human abilities and improve lives. With further enhancements, VisionEcho could become a standard assistive tool for millions of visually impaired users worldwide*