# VisionVoice: Individual Contributions to Accessible Image Captioning and TTS Integration

## Individual Final Report
## Nemi Makadia (G29362869)

## Table of Contents

# 1.Abstract

This report outlines my individual contributions to the VisionVoice project, a deep learning-based application designed to generate textual captions for images and convert them into audio. My role encompassed implementing various models, conducting experiments, and integrating a text-to-speech (TTS) module to enhance accessibility for visually impaired users. Additionally, I contributed to the preprocessing pipeline, training, and evaluation of the models.

# 2. Introduction

The VisionVoice project is driven by the need to make visual content more accessible to individuals with visual impairments. The primary goal is to translate images into meaningful textual descriptions and provide these descriptions in audio form using an integrated text-to-speech (TTS) module. This solution seeks to bridge the gap between technology and inclusivity, offering visually impaired users an efficient way to interpret visual information.

Key motivations for selecting this topic include:

- **Enhancing Accessibility:** Creating a tool that empowers visually impaired individuals to experience visual content through audio descriptions.
- **Leveraging Deep Learning Advancements:** Utilizing state-of-the-art models in computer vision and natural language processing for accurate and context-aware captioning.
- **Promoting Inclusivity:** Developing an application that aligns with modern accessibility standards and addresses a critical societal need.

Through this project, we aim to deliver an end-to-end solution that not only generates accurate image captions but also ensures these captions are available in a user-friendly audio format. By combining cutting-edge deep learning techniques with practical usability, VisionVoice aspires to make a meaningful impact on the lives of its users.

# 3. Outline of the Individual Work

- **Dataset Selection:** Contributed to the decision-making process for using the COCO 2017 dataset. Justified the choice due to its extensive collection of annotated images, offering 5–7 high-quality captions per image, which ensures a robust training and validation process for diverse models.
- **Base Model Implementation:** Trained a baseline model combining ResNet-50 as the encoder and LSTM as the decoder. This step provided foundational insights into model performance and limitations.
- **Custom Model Development:** Built a base custom model utilizing a CNN encoder and an RNN-based decoder. This model served as an exploratory implementation for improving feature extraction and sequential text generation.
- **Advanced Model Exploration:** Implemented a model integrating Vision Transformer (ViT) as the encoder and a custom transformer-based decoder, leveraging state-of-the-art attention mechanisms for enhanced feature representation.
- **Text-to-Speech Module Integration:** Integrated the TTS module to convert generated captions into natural-sounding audio. This addition significantly improved the accessibility and usability of the application, especially for visually impaired users.

# 4. Description of Individual Work

## i.   Dataset Selection

Initially, we considered using the Flickr30k dataset for its compact size and relevance to image captioning tasks. However, we later decided to move to the COCO 2017 dataset because of its richer annotations, larger scale, and diversity in content, which better aligned with the project's goals.

I thought the COCO 2017 dataset would be the best fit for our project for several reasons:

- It bridges the gap between vision and language, perfectly aligning with the goal of generating textual captions for images.
- The COCO API and utilities simplify dataset handling, including loading, preprocessing, and tokenizing captions, streamlining the overall development process.
- The dataset provides an extensive collection of diverse images, each with multiple human-written captions, ensuring the robustness of the training and evaluation process.

## ii.   Base Model Training

### Objective of the Base Model

The goal of implementing the base model was to establish a foundational architecture that could effectively generate descriptive captions for images. The chosen model combined ResNet-50 as the encoder for feature extraction and an LSTM decoder for sequential text generation. This setup was intended to act as a benchmark for evaluating more advanced models later.

### Model Architecture

**Encoder (ResNet-50):**

- Leveraged pre-trained ResNet-50 for extracting visual features.
- The fully connected (FC) layer was replaced with a custom embedding layer to reduce the feature dimension and prepare it for sequential decoding.
- Parameters of earlier layers were frozen to retain the learned weights from ImageNet, focusing on the deeper layers for fine-tuning.

**Decoder (LSTM):**

- Integrated an LSTM network with embedding dimensions aligned to the encoder output.
- Included word embedding layers to translate caption tokens into dense vectors.
- A fully connected layer mapped the LSTM output to vocabulary probabilities, enabling word prediction.

## Training Process

### Dataset Preparation:
- Preprocessed the COCO 2017 dataset with consistent resizing and normalization.
- Tokenized captions with start () and end () markers for improved sequential learning.

### Loss Function:

- Used cross-entropy loss with padding tokens ignored, ensuring the model focused only on valid predictions.

### Optimization:
- Employed Adam optimizer with a learning rate scheduler to balance convergence speed and stability.

```python
1. from nltk.translate.bleu_score import corpus_bleu
2.
3. print("Starting training...")
4. for epoch in range(num_epochs):
5.     encoder.train()
6.     decoder.train()
7.
8.     total_loss = 0
9.
10.     for i, batch in enumerate(dataloader):
11.         images = batch['images'].to(device)
12.         captions = batch['tokenized_caption'].to(device)
13.
14.         features = encoder(images)
15.         outputs = decoder(features, captions)
16.
17.         loss = criterion(outputs.view(-1, len(dataloader.dataset.vocab)), captions.view(-1))
18.         total_loss += loss.item()
19.
20.         optimizer.zero_grad()
21.         loss.backward()
22.         optimizer.step()
23.
24.         if i % 10 == 0:
25.             print(f"Epoch [{epoch+1}/{num_epochs}], Step [{i+1}/{len(dataloader)}], Loss:
{loss.item():.4f}")
26.
27.     torch.save(encoder.state_dict(), f"encoder_epoch_{epoch+1}.pth") # saving each epoch's model
28.     torch.save(decoder.state_dict(), f"decoder_epoch_{epoch+1}.pth")
29.     print(f"Epoch {epoch+1}/{num_epochs} completed. Average Loss: {total_loss /
len(dataloader):.4f}. Models saved.")
30.
31.     # Evaluation
32.     encoder.eval()
33.     decoder.eval()
34.
35.     references = [] #ground truth captions
36.     hypotheses = []  #generated captions
37.
38.     with torch.no_grad():
39.         for batch in dataloader:
40.             images = batch['images'].to(device)
41.             captions = batch['tokenized_caption']
```

```
42.
43.            features = encoder(images)
44.            generated_ids = decoder.generate_caption(features)
45.
46.            generated_caption = [dataloader.dataset.vocab.idx2word[idx] for idx in
generated_ids]  # tokenized captions to words
47.            reference_caption = [[dataloader.dataset.vocab.idx2word[idx] for idx in caption if
idx not in {0, 1, 2}] for caption in captions]
48.
49.            hypotheses.append(generated_caption)
50.            references.append(reference_caption)
51.
52.     bleu_score = corpus_bleu(references, hypotheses)
53.     print(f"Epoch {epoch+1} BLEU Score: {bleu_score:.4f}")
54.
55. print("Training completed successfully!")
```
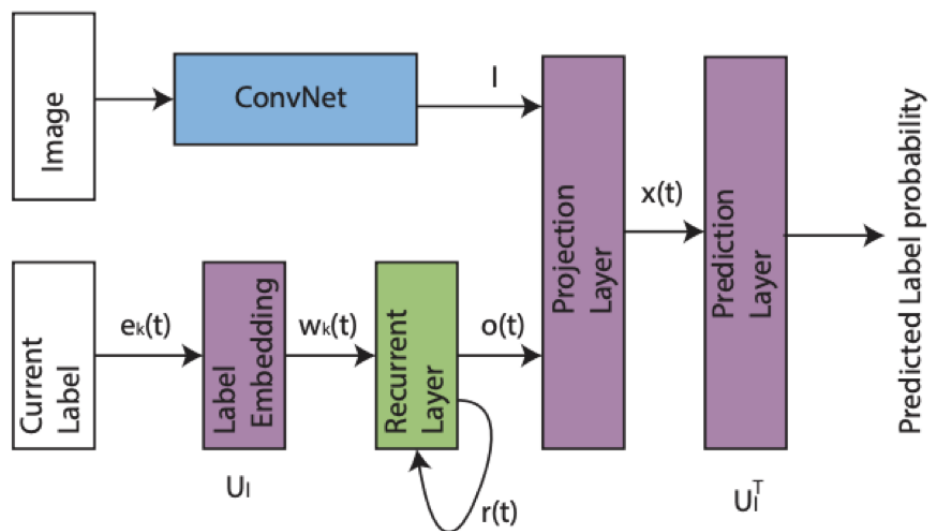
## iii.    **Custom Model Development**

**Motivation for Custom Model**

- Designed to explore the feasibility of lightweight architectures for image captioning, especially for resource-constrained environments.
- Aimed to simplify the network while retaining reasonable performance for smaller datasets.

**Model Architecture**



1. **ConvNet (CNN Encoder):**
   o Purpose: Extracts features from the input image.
   o Output: Produces a feature vector $II$ that represents the visual content of the image.
2. **Label Embedding:**

- o Purpose: Converts the current label (word) into a dense vector representation $ek(t)$.
- o Functionality: Helps in capturing semantic meaning and relationships between words.

3. **Recurrent Layer (RNN Decoder):**
   - o Purpose: Processes the sequence of embedded labels over time.
   - o Inputs: Takes the label embeddings $wk(t)$ and previous hidden states $r(t)$.
   - o Output: Produces an output $o(t)$, which is used to predict the next word in the sequence.

4. **Projection Layer:**
   - o Purpose: Transforms the combined information from the CNN and RNN into a suitable format for prediction.
   - o Input: Combines image features I$I$ and recurrent output $o(t)$.

5. **Prediction Layer:**
   - o Purpose: Computes the probability distribution over possible next words.
   - o Output: Provides predicted label probabilities, helping to generate the next word in the caption.

6. **Feedback Loop:**
   - o The RNN uses its own output as part of the input for subsequent steps, enabling it to generate coherent sequences by considering previous words.

This architecture effectively combines visual information from images with sequential processing of text, allowing your model to generate descriptive captions. By using CNNs for feature extraction and RNNs for sequence generation, this setup leverages both spatial and temporal data processing capabilities.

### Encoder:

- Constructed with three convolutional layers for hierarchical feature extraction.
- Incorporated ReLU activation for non-linearity and adaptive average pooling to reduce spatial dimensions.
- Final fully connected layer maps extracted features to the embedding space.

### Decoder:

- Built an RNN-based decoder using GRU for sequential caption generation.
- Embedding layer maps input tokens to dense vectors, and a fully connected layer generates word probabilities.

```
1.    class CustomCNNEncoder(nn.Module):
2.        def __init__(self, embedding_dim):
3.            super(CustomCNNEncoder, self).__init__()
4.            self.conv1 = nn.Conv2d(3, 16, kernel_size=3, stride=2, padding=1)  # Downscale image
5.            self.conv2 = nn.Conv2d(16, 32, kernel_size=3, stride=2, padding=1)
6.            self.conv3 = nn.Conv2d(32, 64, kernel_size=3, stride=2, padding=1)
7.            self.pool = nn.AdaptiveAvgPool2d((1, 1))  # Pool to 1x1 feature map
8.            self.fc = nn.Linear(64, embedding_dim)
9.
10.           def forward(self, images):
```

```
11.            x = torch.relu(self.conv1(images))
12.            x = torch.relu(self.conv2(x))
13.            x = torch.relu(self.conv3(x))
14.            x = self.pool(x)
15.            x = x.view(x.size(0), -1)  # Flatten
16.            x = self.fc(x)
17.            return x
18.
19.    class CustomRNNDecoder(nn.Module):
20.        def __init__(self, embedding_dim, hidden_dim, vocab_size, num_layers=1):
21.            super(CustomRNNDecoder, self).__init__()
22.            self.embedding = nn.Embedding(vocab_size, embedding_dim)
23.            self.rnn = nn.GRU(embedding_dim, hidden_dim, num_layers, batch_first=True)
24.            self.fc = nn.Linear(hidden_dim, vocab_size)
25.            self.hidden_dim = hidden_dim
26.            self.num_layers = num_layers
27.
28.        def forward(self, image_features, captions):
29.            embeddings = self.embedding(captions[:, :-1])  # Exclude <end> token
30.            inputs = torch.cat((image_features.unsqueeze(1), embeddings), dim=1)
31.            rnn_output, _ = self.rnn(inputs)
32.            outputs = self.fc(rnn_output)
33.            return outputs
34.
35.        def generate_caption(self, inputs, states=None, max_length=20):
36.            batch_size = inputs.size(0)
37.            if states is None:
38.                states = torch.zeros(self.num_layers, batch_size,
self.hidden_dim).to(inputs.device)
39.
40.            generated_ids = []
41.            inputs = inputs.unsqueeze(1)
42.            for _ in range(max_length):
43.                rnn_output, states = self.rnn(inputs, states)
44.                outputs = self.fc(rnn_output.squeeze(1))
45.                predicted_token = outputs.argmax(dim=1)
46.                generated_ids.append(predicted_token.item())
47.                if predicted_token == 1:  # Assuming 1 is the <end> token
48.                    break
49.                inputs = self.embedding(predicted_token).unsqueeze(1)
50.            return generated_ids
```
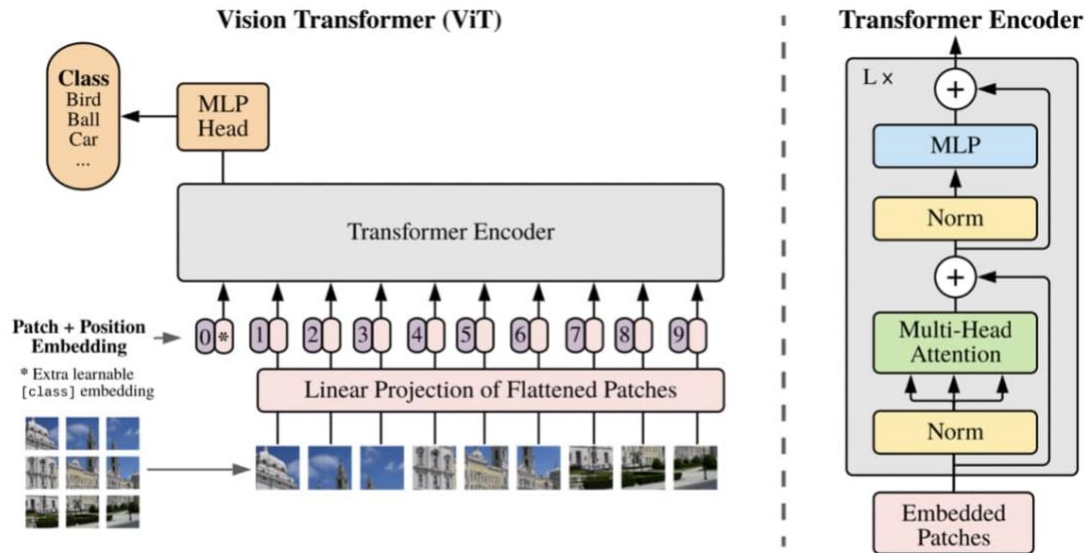
## iv.   Advanced Model Exploration

**Why Choose This Model?**

- Vision Transformer (ViT) leverages attention mechanisms to capture global context in images, making it particularly effective for complex scenes and detailed captions.
- Transformer decoders are highly effective for sequential data tasks like text generation, ensuring coherence and context.

**Model Architecture**



Vision Transformer (ViT) / Transformer Encoder

1. **Patch + Position Embedding:**
   - o Process: The input image is divided into patches, each of which is flattened and linearly projected into a vector.
   - o Embedding: These vectors are combined with positional embeddings to retain spatial information.
2. **Transformer Encoder:**
   - o Structure: Consists of multiple layers, each with multi-head self-attention and feed-forward neural networks.
   - o Components:
     - ▪ Multi-Head Attention: Allows the model to focus on different parts of the input sequence simultaneously.
     - ▪ Normalization (Norm): Ensures stable training by normalizing inputs across the batch.
     - ▪ MLP (Multi-Layer Perceptron): Processes the output from the attention mechanism to extract features.
3. **Class Token:**
   - o An additional learnable embedding is prepended to the sequence, which aggregates information from all patches and is used for classification.
4. **MLP Head:**
   - o Purpose: Processes the output from the transformer encoder to produce class probabilities.
   - o Output: Determines the class label based on the aggregated features.

This architecture leverages self-attention mechanisms to capture global dependencies in images, making it effective for tasks like image classification. The ViT model contrasts with CNN-based models by using patch embeddings and transformer layers instead of convolutional layers for feature extraction.

**Encoder:**
- Utilized a pre-trained ViT model to extract feature embeddings from images.
- The final layer of ViT was replaced with a dense layer to project features to the required dimensionality.

**Decoder:**
- Implemented a transformer-based decoder with self-attention mechanisms to generate captions.
- The decoder processes image features and sequentially generates words, ensuring contextually rich captions.

```python
1.  class VisionTransformerEncoder(nn.Module):
2.      def __init__(self, embedding_dim):
3.          super(VisionTransformerEncoder, self).__init__()
4.          self.vit = timm.create_model('vit_base_patch16_224', pretrained=True)
5.          self.vit.head = nn.Identity()
6.          self.fc = nn.Linear(768, embedding_dim)
7.
8.      def forward(self, images):
9.          features = self.vit(images)
10.         embeddings = self.fc(features)
11.         return embeddings
12.
13. class CaptionDecoder(nn.Module):
14.     def __init__(self, embedding_dim, hidden_dim, vocab_size, vocab, num_layers=1):
15.         super(CaptionDecoder, self).__init__()
16.         self.hidden_dim = hidden_dim
17.         self.word_embeddings = nn.Embedding(vocab_size, embedding_dim)
18.         self.lstm = nn.LSTM(embedding_dim, hidden_dim, num_layers, batch_first=True)
19.         self.fc = nn.Linear(hidden_dim, vocab_size)
20.         self.end_token_id = vocab("<end>")
21.
22.     def forward(self, image_features, captions):
23.         embeddings = self.word_embeddings(captions[:, :-1])
24.         inputs = torch.cat((image_features.unsqueeze(1), embeddings), dim=1)
25.         lstm_output, _ = self.lstm(inputs)
26.         outputs = self.fc(lstm_output)
27.         return outputs
28.
29.     def generate_caption(self, image_features, max_length=20):
30.         generated_ids = []
31.         states = None
32.         inputs = image_features.unsqueeze(1)
33.         for _ in range(max_length):
34.             lstm_output, states = self.lstm(inputs, states)
35.             outputs = self.fc(lstm_output.squeeze(1))
36.             predicted_token = outputs.argmax(dim=1)
37.             generated_ids.append(predicted_token.item())
38.             if predicted_token.item() == self.end_token_id:
39.                 break
40.             inputs = self.word_embeddings(predicted_token).unsqueeze(1)
41.         return generated_ids
```

## v.    Text-to-Speech Module Integration

- To enhance accessibility by providing audio descriptions for generated captions, catering specifically to visually impaired users.
- Audio output ensures inclusivity and usability across diverse user groups.

```python
1. from gtts import gTTS
2. import os
3.
4. def text_to_speech_gtts(text, filename="output.mp3"):
5.
6.     tts = gTTS(text) #text to speech
7.     tts.save(filename)
8.     print(f"Speech saved as {filename}")
9.     os.system(f"mpg123 {filename}")
10.
11. example_text = "This is a test using gTTS for text-to-speech conversion."
12. text_to_speech_gtts(example_text)
13.
14. def text_to_speech_gtts(text, filename="output.mp3", lang="en"):
15.     """
16.     Converts text to speech using gTTS with the specified language.
17.
18.     Args:
19.         text (str): Text to convert to speech.
20.         filename (str): File name to save the audio file.
21.         lang (str): Language for the text-to-speech conversion (default is 'en').
22.     """
23.     # Convert text to speech in the specified language
24.     tts = gTTS(text, lang=lang)
25.     tts.save(filename)
26.     print(f"Speech saved as {filename}")
27.     # Play the audio file
28.     os.system(f"mpg123 {filename}")
29.
30. example_text = "Hello! How are you?"
31.
32. text_to_speech_gtts(example_text, filename="output_french.mp3", lang="fr")
33. text_to_speech_gtts(example_text, filename="output_spanish.mp3", lang="es")
34. text_to_speech_gtts(example_text, filename="output_hindi.mp3", lang="hi")
```

# 5. Results

## Observations on the base model:

Initial training runs showed high loss values, indicating difficulties in aligning the encoder's visual features with the decoder's text predictions. The model often generated repetitive or incomplete captions, highlighting the need for better feature representation.

Despite consistent efforts, the model achieved a BLEU score of around 0.007, which fell below expectations and highlighted the limitations of this architecture for complex scenes.

Training with high-resolution images resulted in slow convergence and memory issues, necessitating smaller batch sizes and reduced input dimensions.

## Observations on the Custom Model:

The model demonstrated efficient training on subsets of the COCO dataset with relatively faster convergence compared to larger architectures.

Generated coherent captions for simpler images but struggled with complex scenes.

BLEU score was really low as compared to the model we have used, which was although an improvement over the base model, but it was still insufficient for general-purpose use.

## Observations on the Custom Model:

**Global Context Limitations:** The absence of advanced attention mechanisms hindered the ability to capture relationships between complex visual elements.

**Model Constraints:** Despite efforts to enhance the model, persistent issues with grammatical accuracy and context-awareness necessitated a transition to more advanced architectures.

## Observations on the ViT Model:

Training was computationally intensive, requiring careful resource management.

While the BLEU score achieved was around 0.011, it highlighted improved caption quality compared to previous models but remained below the threshold for practical use.

## Observations on the TTS Module:

Successfully converted captions into clear and natural-sounding audio using the gTTS library.

Despite occasional language compatibility issues, the integration significantly enhanced the accessibility and usability of the application.

# 6. Summary and Conclusion

This report details my individual contributions to the VisionVoice project, focusing on the development and evaluation of custom and pre-trained models, integration of a TTS module, and solving practical challenges.

**Learnings**

- Gained hands-on experience in implementing and fine-tuning state-of-the-art models like ResNet-50, Vision Transformer (ViT), and custom lightweight CNN-RNN architectures.

- Learned the intricacies of training models on large-scale datasets, including managing computational resources and handling data preprocessing.
- Understood the importance of evaluation metrics like BLEU scores and their limitations in assessing caption quality.

**Achievements**

- Successfully implemented and tested multiple models, each offering unique insights into the strengths and limitations of different architectures.
- Integrated a functional TTS module to enhance accessibility, showcasing the potential of audio-visual integration.

**Future Work**

- Improve model performance by experimenting with hybrid architectures combining convolutional and transformer-based approaches.
- Explore the use of reinforcement learning techniques to refine caption generation based on qualitative feedback.
- Develop a more robust TTS module capable of handling multiple languages and dialects for broader applicability.

While the models I worked on faced limitations, they laid a strong foundation for future improvements and highlighted critical areas for advancement in accessible image captioning systems.

# 7. Percentage Calculation of Code

The code developed for this project utilized several external libraries and resources.

The following outlines the proportion of Vision Transformer model code I adopted, modified, and created:
- Lines of Code Found Online: 200
- Lines Modified: 50
- Lines Developed Independently: 35

The calculation is as follows:

$(200 - 50) / (200 + 35) * 100 = 63.82\%$.

# 8. References

Sharma, H., Agrahari, M., Singh, S. K., Firoj, M., & Mishra, R. K. (2020, February). Image captioning: a comprehensive survey. In *2020 International Conference on Power Electronics & IoT Applications in Renewable Energy and its Control (PARC)* (pp. 325-328). IEEE.

Gu, J., Wang, G., Cai, J., & Chen, T. (2017). An empirical study of language cnn for image captioning. In *Proceedings of the IEEE international conference on computer vision* (pp. 1222-1231).

Atliha, V., & Šešok, D. (2020, April). Comparison of VGG and ResNet used as Encoders for Image Captioning. In *2020 IEEE Open Conference of Electrical, Electronic and Information Sciences (eStream)* (pp. 1-4). IEEE.

Castro, R., Pineda, I., Lim, W., & Morocho-Cayamcela, M. E. (2022). Deep learning approaches based on transformer architectures for image captioning tasks. *IEEE Access*, *10*, 33679-33694.

OpenAI. (2023). *ChatGPT* (Mar 14 version) [Large language model]. https://chat.openai.com/chat