

EXP NO: 07  
DATE: 20/09/2025

## CAPTION GENERATION USING RNN+CNN

**AIM:** To build a deep recurrent neural network (RNN) that generates image captions by combining a Convolutional Neural Network (CNN) for image feature extraction with an RNN for sequence generation using the MS COCO (or Flickr8k) dataset.

### ALGORITHM:

- Import TensorFlow, Keras, and supporting libraries.
- Load and preprocess the caption dataset, adding start and end tokens to each caption.
- Use a pre-trained CNN model (InceptionV3) to extract feature vectors from images.
- Tokenize and pad captions, converting text to numerical sequences.
- Create input-output pairs combining image features and partial text sequences for training.
- Define a multimodal model combining CNN features and RNN outputs: CNN output passes through a dense layer, Captions are embedded and processed through an LSTM, their outputs are merged and passed to a softmax layer for word prediction.
- Train the model using sparse categorical cross-entropy loss.
- Generate captions by iteratively predicting the next word until the end token is reached.
- Display the image with the generated caption for evaluation.

### CODE:

```
import tensorflow as tf

from tensorflow.keras.applications import InceptionV3

from tensorflow.keras.applications.inception_v3 import preprocess_input

from tensorflow.keras.preprocessing.text import Tokenizer

from tensorflow.keras.preprocessing.sequence import pad_sequences

from tensorflow.keras.layers import Input, Embedding, LSTM, Dense, Add

from tensorflow.keras.models import Model

import numpy as np, os, pandas as pd

from PIL import Image

from tqdm import tqdm

import matplotlib.pyplot as plt
```

```
# GPU configuration

gpus = tf.config.list_physical_devices('GPU')

if gpus:
    try:
        for gpu in gpus:
            tf.config.experimental.set_memory_growth(gpu, True)
            tf.keras.mixed_precision.set_global_policy('mixed_float16')
    except RuntimeError as e:
        print(f"GPU configuration error: {e}")

print(f"TensorFlow version: {tf.__version__}")

# Paths
IMG_DIR = "Images"
CAP_FILE = "captions.txt"

# Load captions
df = pd.read_csv(CAP_FILE)
df['caption'] = df['caption'].apply(lambda x: 'startseq ' + x.lower() + ' endseq')

# Subset for demo
df = df.groupby('image').head(1).sample(2000, random_state=42)

# CNN feature extractor
cnn = InceptionV3(weights='imagenet', include_top=False, pooling='avg')

def extract_feat(img_path):
```

30

```
img = Image.open(img_path).convert('RGB').resize((299,299))

x = np.expand_dims(preprocess_input(np.array(img)), 0)

return cnn.predict(x, verbose=0)[0]


# Extract features
features, captions = [], []
for img, cap in tqdm(zip(df['image'], df['caption']), total=len(df)):
    path = os.path.join(IMG_DIR, img)
    if os.path.exists(path):
        features.append(extract_feat(path))
        captions.append(cap)
features = np.array(features)


# Tokenize captions
tok = Tokenizer(num_words=5000, oov_token='unk')
tok.fit_on_texts(captions)
seqs = tok.texts_to_sequences(captions)
maxlen = max(len(s) for s in seqs)
vocab = len(tok.word_index) + 1


# Prepare training data
X1, X2, y = [], [], []
for f, s in zip(features, seqs):
    for i in range(1, len(s)):
        in_seq, out = s[:i], s[i]
        X1.append(f)
        X2.append(pad_sequences([in_seq], maxlen=maxlen)[0])
        y.append(out)
```

```
X1, X2, y = np.array(X1), np.array(X2), np.array(y)

# Define CNN+RNN model

img_in = Input(shape=(2048,))
cap_in = Input(shape=(maxlen,))
emb = Embedding(vocab, 256, mask_zero=True)(cap_in)
lstm = LSTM(256)(emb)
x = Add()([Dense(256, activation='relu')(img_in), lstm])
out = Dense(vocab, activation='softmax', dtype='float32')(x)
model = Model([img_in, cap_in], out)
model.compile(loss='sparse_categorical_crossentropy', optimizer='adam')
model.summary()

# Train model

model.fit([X1, X2], y, batch_size=128, epochs=5, verbose=1)

# Caption generator

def generate_caption(img_path):
    f = extract_feat(img_path).reshape(1,-1)
    cap = ['startseq']
    for _ in range(maxlen):
        seq = pad_sequences([tok.texts_to_sequences([' '.join(cap)])][0], maxlen=maxlen)
        pred = np.argmax(model.predict([f, seq], verbose=0))
        word = tok.index_word.get(pred, '')
        cap.append(word)
        if word == 'endseq': break
    return ' '.join(cap[1:-1])
```

# Test

```
test_img = os.path.join(IMG_DIR, df.iloc[0]['image'])
```

```
caption = generate_caption(test_img)
```

```
print("Generated caption:", caption)
```

# Display image with caption

```
img = Image.open(test_img)
```

```
plt.figure(figsize=(10, 6))
```

```
plt.imshow(img)
```

```
plt.axis('off')
```

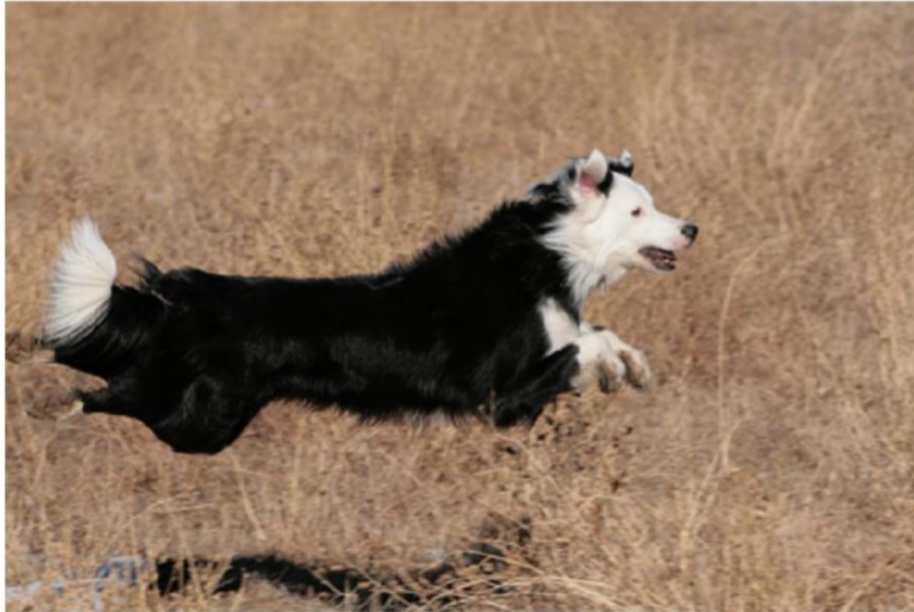
```
plt.title(f"Generated Caption:\n{caption}", fontsize=14, weight='bold', pad=20)
```

```
plt.tight_layout()
```

```
plt.show()
```

**OUTPUT:**

**Generated Caption:**  
**a black and white dog is jumping in the air**



**RESULT:** The CNN+RNN model successfully generated meaningful captions for images, demonstrating the effectiveness of combining CNN for feature extraction and RNN for sequence generation.