



[5]:

```
# This Python 3 environment comes with many helpful analytics libraries installed
# It is defined by the kaggle/python Docker image: https://github.com/kaggle/docker-python
# For example, here's several helpful packages to load

import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)

# Input data files are available in the read-only "../input/" directory
# For example, running this (by clicking run or pressing Shift+Enter) will list all files under the input directory

import os
# for dirname, _, filenames in os.walk('/kaggle/input'):
#     # for filename in filenames:
#         # print(os.path.join(dirname, filename))

# You can write up to 20GB to the current directory (/kaggle/working/) that gets preserved as output when you create a version using "Save"
# You can also write temporary files to /kaggle/temp/, but they won't be saved outside of the current session
```

[6]:

```
import numpy as np
import pandas as pd
import os
import tensorflow as tf
from tqdm import tqdm
from tensorflow.keras.preprocessing.image import ImageDataGenerator, load_img, img_to_array
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.utils import Sequence
from tensorflow.keras.utils import to_categorical
from tensorflow.keras.models import Sequential, Model
from tensorflow.keras.layers import Conv2D, MaxPooling2D, GlobalAveragePooling2D, Activation, Dropout, Flatten, Dense, Input, Layer
from tensorflow.keras.layers import Embedding, LSTM, add, Concatenate, Reshape, concatenate, Bidirectional
from tensorflow.keras.applications import VGG16, ResNet50, DenseNet201
from tensorflow.keras.optimizers import Adam
from tensorflow.keras.callbacks import ModelCheckpoint, EarlyStopping, ReduceLROnPlateau
import warnings
import matplotlib.pyplot as plt
import seaborn as sns
from textwrap import wrap

plt.rcParams['font.size'] = 12
sns.set_style("dark")
warnings.filterwarnings('ignore')
```

```
/opt/conda/lib/python3.10/site-packages/tensorflow_io/python/ops/_init__.py:98: UserWarning: unable to load libtensorflow_io_plugins.so: unable to open file: libtensorflow_io_plugins.so, from paths: ['/opt/conda/lib/python3.10/site-packages/tensorflow_io/python/ops/libtensorflow_io_plugins.so']
caused by: ['/opt/conda/lib/python3.10/site-packages/tensorflow_io/python/ops/libtensorflow_io_plugins.so: undefined symbol: _ZN3tsl6StatusC1EN10tensorflowError4CodeSt17basic_string_viewIcSt1char_traitsIcEENS_14SourceLocationE']
  warnings.warn(f"unable to load libtensorflow_io_plugins.so: {e}")
/opt/conda/lib/python3.10/site-packages/tensorflow_io/python/ops/_init__.py:104: UserWarning: file system plugins are not loaded: unable to open file: libtensorflow_io.so, from paths: ['/opt/conda/lib/python3.10/site-packages/tensorflow_io/python/ops/libtensorflow_io.so']
caused by: ['/opt/conda/lib/python3.10/site-packages/tensorflow_io/python/ops/libtensorflow_io.so: undefined symbol: _ZTVN10tensorflow13GcsFileSystemE']
  warnings.warn(f"file system plugins are not loaded: {e}")
```

[7]:

```
image_path = '../input/flickr8k/Images'
```

[8]:

```
data = pd.read_csv("../input/flickr8k/captions.txt")
data.head()
```

[8]:

	image	caption
0	1000268201_693b08cb0e.jpg	A child in a pink dress is climbing up a set o...
1	1000268201_693b08cb0e.jpg	A girl going into a wooden building .
2	1000268201_693b08cb0e.jpg	A little girl climbing into a wooden playhouse .
3	1000268201_693b08cb0e.jpg	A little girl climbing the stairs to her play...
4	1000268201_693b08cb0e.jpg	A little girl in a pink dress going into a woo...

[9]:

```
def readImage(path,img_size=224):
    img = load_img(path,color_mode='rgb',target_size=(img_size,img_size))
    img = img_to_array(img)
    img = img/255.
```

```

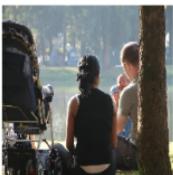
return img

def display_images(temp_df):
    temp_df = temp_df.reset_index(drop=True)
    plt.figure(figsize = (20 , 20))
    n = 0
    for i in range(15):
        n+=1
        plt.subplot(5 , 5, n)
        plt.subplots_adjust(hspace = 0.7, wspace = 0.3)
        image = readImage(f"../input/flickr8k/Images/{temp_df.image[i]}")
        plt.imshow(image)
        plt.title("\n".join(wrap(temp_df.caption[i], 20)))
        plt.axis("off")

```

[10]: display\_images(data.sample(15))

Couple with a baby sit outdoors next to their stroller .



A dog swims in the water .



A hockey game is being played .



A protester in a cowboy gear holds up a sign that says , "Dont tax me bro!"



Two young boys read comic books in bed .



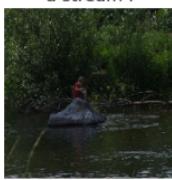
a boy rides his ripstik on the street .



A boy doing a handstand in a yellow shirt



A child on a rock in a stream .



A man cleaning the inside of a large pane of glass .



Two women ride camels , while a third watches .



A young girl sitting in swimming gear



Boy sliding down a slide on his knees .



a lone bicyclist jumping a trick along a bike jumping path .



A woman stands against a wall near a video camera .



A man lays on a bench while his dog sits by him .



[11]:

```

def text_preprocessing(data):
    data['caption'] = data['caption'].apply(lambda x: x.lower())
    data['caption'] = data['caption'].apply(lambda x: x.replace("[^A-Za-z]", ""))
    data['caption'] = data['caption'].apply(lambda x: x.replace("\s+", " "))
    data['caption'] = data['caption'].apply(lambda x: " ".join([word for word in x.split() if len(word)>1]))
    data['caption'] = "startseq "+data['caption']+ " endseq"
    return data

```

## Preprocessed Text

[12]:

```

data = text_preprocessing(data)
captions = data['caption'].tolist()
captions[:10]

```

[12]:

```

'startseq child in pink dress is climbing up set of stairs in an entry way endseq',
'startseq girl going into wooden building endseq',
'startseq little girl climbing into wooden playhouse endseq',
'startseq little girl climbing the stairs to her playhouse endseq',
'startseq little girl in pink dress going into wooden cabin endseq',
'startseq black dog and spotted dog are fighting endseq',
'startseq black dog and tri-colored dog playing with each other on the road endseq',
'startseq black dog and white dog with brown spots are staring at each other in the street endseq',
'startseq two dogs of different breeds looking at each other on the road endseq',
'<startseq two dogs on pavement moving toward each other endseq'

```

```
[13]: 
tokenizer = Tokenizer()
tokenizer.fit_on_texts(captions)
vocab_size = len(tokenizer.word_index) + 1
max_length = max(len(caption.split()) for caption in captions)

images = data['image'].unique().tolist()
nimages = len(images)

split_index = round(0.85*nimages)
train_images = images[:split_index]
val_images = images[split_index:]

train = data[data['image'].isin(train_images)]
test = data[data['image'].isin(val_images)]

train.reset_index(inplace=True, drop=True)
test.reset_index(inplace=True, drop=True)

tokenizer.texts_to_sequences([captions[1]])[0]
```

[13]: [1, 18, 315, 63, 195, 116, 2]

```
[14]: 
model = DenseNet201()
fe = Model(inputs=model.input, outputs=model.layers[-2].output)

img_size = 224
features = {}
for image in tqdm(data['image'].unique().tolist()):
    img = load_img(os.path.join(image_path, image), target_size=(img_size, img_size))
    img = img_to_array(img)
    img = img/255.
    img = np.expand_dims(img, axis=0)
    feature = fe.predict(img, verbose=0)
    features[image] = feature
```

Downloading data from [https://storage.googleapis.com/tensorflow/keras-applications/densenet/densenet201\\_weights\\_tf\\_dim\\_ordering\\_tf\\_kernels.h5](https://storage.googleapis.com/tensorflow/keras-applications/densenet/densenet201_weights_tf_dim_ordering_tf_kernels.h5)
82524592/82524592 [=====] - 1s 0us/step
100%|██████████| 8091/8091 [13:56<00:00, 9.67it/s]

```
[15]: 
class CustomDataGenerator(Sequence):

    def __init__(self, df, X_col, y_col, batch_size, directory, tokenizer,
                 vocab_size, max_length, features, shuffle=True):

        self.df = df.copy()
        self.X_col = X_col
        self.y_col = y_col
        self.directory = directory
        self.batch_size = batch_size
        self.tokenizer = tokenizer
        self.vocab_size = vocab_size
        self.max_length = max_length
        self.features = features
        self.shuffle = shuffle
        self.n = len(self.df)

    def on_epoch_end(self):
        if self.shuffle:
            self.df = self.df.sample(frac=1).reset_index(drop=True)

    def __len__(self):
        return self.n // self.batch_size

    def __getitem__(self, index):

        batch = self.df.iloc[index * self.batch_size:(index + 1) * self.batch_size, :]
        X1, X2, y = self.__get_data(batch)
        return (X1, X2), y

    def __get_data(self, batch):

        X1, X2, y = list(), list(), list()

        images = batch[self.X_col].tolist()

        for image in images:
            feature = self.features[image][0]

            captions = batch.loc[batch[self.X_col]==image, self.y_col].tolist()
            for caption in captions:
                seq = self.tokenizer.texts_to_sequences([caption])[0]
```

```

        for i in range(1,len(seq)):
            in_seq, out_seq = seq[:i], seq[i]
            in_seq = pad_sequences([in_seq], maxlen=self.max_length)[0]
            out_seq = to_categorical([out_seq], num_classes=self.vocab_size)[0]
            X1.append(feature)
            X2.append(in_seq)
            y.append(out_seq)

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

    return X1, X2, y

```

## Modelling

- 1.The image embedding representations are concatenated with the first word of sentence ie. starseq and passed to the LSTM network
- 2.The LSTM network starts generating words after each input thus forming a sentence at the end

```
[16]:
input1 = Input(shape=(1920,))
input2 = Input(shape=(max_length,))

img_features = Dense(256, activation='relu')(input1)
img_features_reshaped = Reshape((1, 256), input_shape=(256,))(img_features)

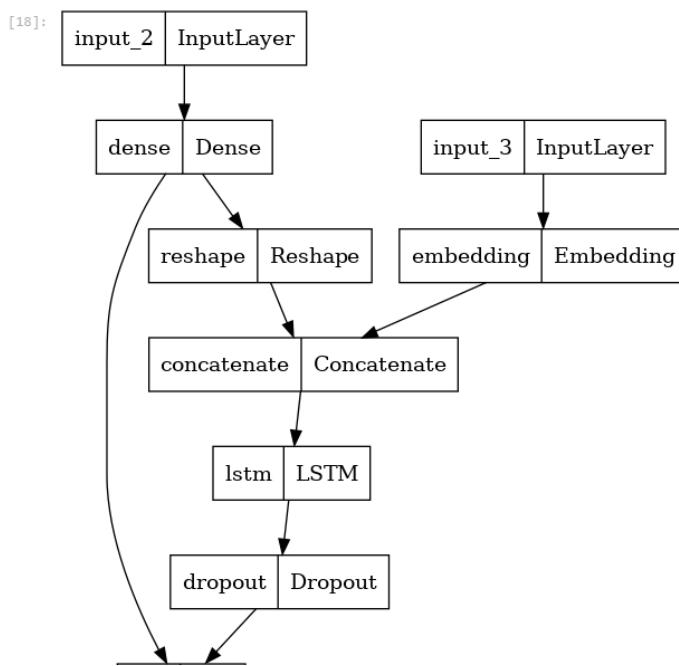
sentence_features = Embedding(vocab_size, 256, mask_zero=False)(input2)
merged = concatenate([img_features_reshaped, sentence_features], axis=1)
sentence_features = LSTM(256)(merged)
x = Dropout(0.5)(sentence_features)
x = add([x, img_features])
x = Dense(128, activation='relu')(x)
x = Dropout(0.5)(x)
output = Dense(vocab_size, activation='softmax')(x)

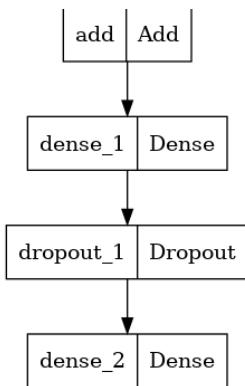
caption_model = Model(inputs=[input1,input2], outputs=output)
caption_model.compile(loss='categorical_crossentropy',optimizer='adam')
```

```
[17]:
from tensorflow.keras.utils import plot_model
```

## Model Modification

```
[18]:
plot_model(caption_model)
```





```
[19]: caption_model.summary()
```

```
Model: "model_1"
-----
```

Layer (type)	Output Shape	Param #	Connected to
input_2 (InputLayer)	(None, 1920)	0	[]
dense (Dense)	(None, 256)	491776	['input_2[0][0]']
input_3 (InputLayer)	(None, 34)	0	[]
reshape (Reshape)	(None, 1, 256)	0	['dense[0][0]']
embedding (Embedding)	(None, 34, 256)	2172160	['input_3[0][0]']
concatenate (Concatenate)	(None, 35, 256)	0	['reshape[0][0]', 'embedding[0][0]']
lstm (LSTM)	(None, 256)	525312	['concatenate[0][0]']
dropout (Dropout)	(None, 256)	0	['lstm[0][0]']
add (Add)	(None, 256)	0	['dropout[0][0]', 'dense[0][0]']
dense_1 (Dense)	(None, 128)	32896	['add[0][0]']
dropout_1 (Dropout)	(None, 128)	0	['dense_1[0][0]']
dense_2 (Dense)	(None, 8485)	1094565	['dropout_1[0][0]']

```
=====
Total params: 4,316,709
Trainable params: 4,316,709
Non-trainable params: 0
```

```
[20]:
train_generator = CustomDataGenerator(df=train,X_col='image',y_col='caption',batch_size=64,directory=image_path,
                                      tokenizer=tokenizer,vocab_size=vocab_size,max_length=max_length,features=features)

validation_generator = CustomDataGenerator(df=test,X_col='image',y_col='caption',batch_size=64,directory=image_path,
                                           tokenizer=tokenizer,vocab_size=vocab_size,max_length=max_length,features=features)
```

```
[21]:
model_name = "model.h5"
checkpoint = ModelCheckpoint(model_name,
                            monitor='val_loss',
                            mode="min",
                            save_best_only = True,
                            verbose=1)

earlystopping = EarlyStopping(monitor='val_loss',min_delta = 0, patience = 5, verbose = 1, restore_best_weights=True)

learning_rate_reduction = ReduceLROnPlateau(monitor='val_loss',
                                             patience=3,
                                             verbose=1,
                                             factor=0.2,
                                             min_lr=0.0000001)
```

```
[22]:
history = caption_model.fit(
    train_generator,
    epochs=50,
    validation_data=validation_generator,
    callbacks=[checkpoint,earlystopping,learning_rate_reduction])
```

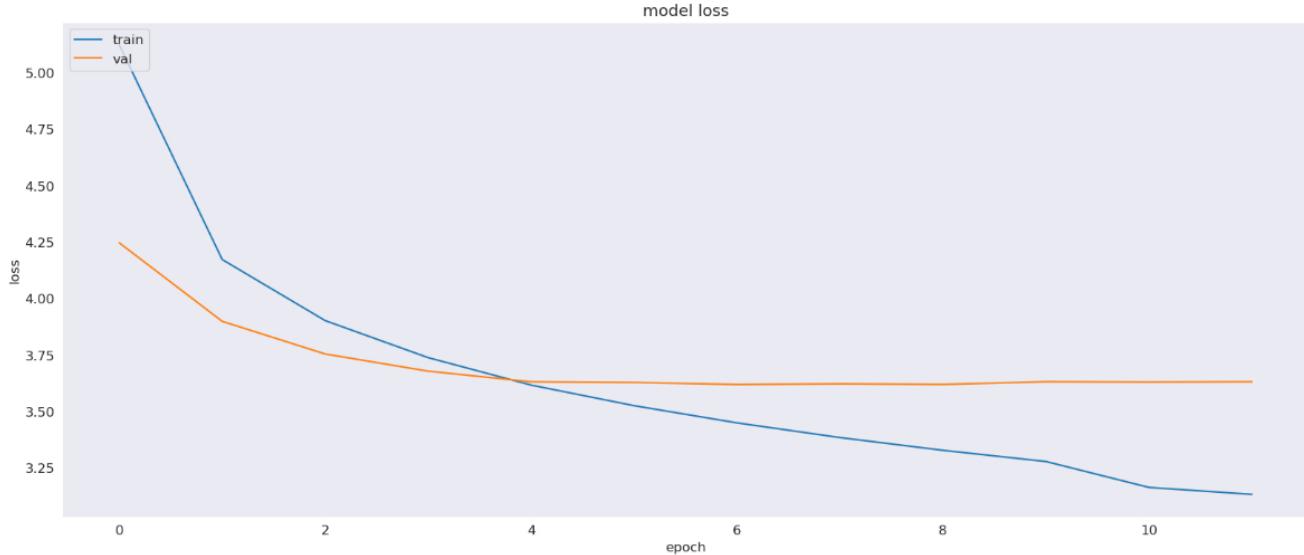
```
Epoch 1/50
537/537 [=====] - ETA: 0s - loss: 5.1236
Epoch 1: val_loss improved from inf to 4.24643, saving model to model.h5
537/537 [=====] - 255s 466ms/step - loss: 5.1236 - val_loss: 4.2464 - lr: 0.00010
Epoch 2/50
```

```

537/537 [=====] - ETA: 0s - loss: 4.1722
Epoch 2: val_loss improved from 4.24643 to 3.89864, saving model to model.h5
537/537 [=====] - 68s 127ms/step - loss: 4.1722 - val_loss: 3.8986 - lr: 0.0010
Epoch 3/50
537/537 [=====] - ETA: 0s - loss: 3.9018
Epoch 3: val_loss improved from 3.89864 to 3.75393, saving model to model.h5
537/537 [=====] - 51s 96ms/step - loss: 3.9018 - val_loss: 3.7539 - lr: 0.0010
Epoch 4/50
537/537 [=====] - ETA: 0s - loss: 3.7378
Epoch 4: val_loss improved from 3.75393 to 3.67802, saving model to model.h5
537/537 [=====] - 51s 96ms/step - loss: 3.7378 - val_loss: 3.6780 - lr: 0.0010
Epoch 5/50
537/537 [=====] - ETA: 0s - loss: 3.6162
Epoch 5: val_loss improved from 3.67802 to 3.63076, saving model to model.h5
537/537 [=====] - 50s 93ms/step - loss: 3.6162 - val_loss: 3.6308 - lr: 0.0010
Epoch 6/50
537/537 [=====] - ETA: 0s - loss: 3.5250
Epoch 6: val_loss improved from 3.63076 to 3.62794, saving model to model.h5
537/537 [=====] - 51s 95ms/step - loss: 3.5250 - val_loss: 3.6279 - lr: 0.0010
Epoch 7/50
537/537 [=====] - ETA: 0s - loss: 3.4488
Epoch 7: val_loss improved from 3.62794 to 3.61881, saving model to model.h5
537/537 [=====] - 50s 94ms/step - loss: 3.4488 - val_loss: 3.6188 - lr: 0.0010
Epoch 8/50
537/537 [=====] - ETA: 0s - loss: 3.3838
Epoch 8: val_loss did not improve from 3.61881
537/537 [=====] - 50s 93ms/step - loss: 3.3838 - val_loss: 3.6218 - lr: 0.0010
Epoch 9/50
537/537 [=====] - ETA: 0s - loss: 3.3273
Epoch 9: val_loss did not improve from 3.61881
537/537 [=====] - 51s 95ms/step - loss: 3.3273 - val_loss: 3.6191 - lr: 0.0010
Epoch 10/50
537/537 [=====] - ETA: 0s - loss: 3.2777
Epoch 10: val_loss did not improve from 3.61881
Epoch 10: ReducelROnPlateau reducing learning rate to 0.0002000000949949026.
537/537 [=====] - 50s 92ms/step - loss: 3.2777 - val_loss: 3.6312 - lr: 0.0010
Epoch 11/50
537/537 [=====] - ETA: 0s - loss: 3.1630
Epoch 11: val_loss did not improve from 3.61881
537/537 [=====] - 50s 92ms/step - loss: 3.1630 - val_loss: 3.6298 - lr: 2.0000e-04
Epoch 12/50
537/537 [=====] - ETA: 0s - loss: 3.1323
Epoch 12: val_loss did not improve from 3.61881
Restoring model weights from the end of the best epoch: 7.
537/537 [=====] - 50s 93ms/step - loss: 3.1323 - val_loss: 3.6311 - lr: 2.0000e-04
Epoch 12: early stopping

```

```
[23]: plt.figure(figsize=(20,8))
plt.plot(history.history['loss'])
plt.plot(history.history['val_loss'])
plt.title('model loss')
plt.ylabel('loss')
plt.xlabel('epoch')
plt.legend(['train', 'val'], loc='upper left')
plt.show()
```



```
[24]: def idx_to_word(integer,tokenizer):
    for word, index in tokenizer.word_index.items():
        if index==integer:
            return word
    return None
```

```
[25]: def predict_caption(model, image, tokenizer, max_length, features):
    feature = features[image]
```

```
in_text = startseq
for i in range(max_length):
    sequence = tokenizer.texts_to_sequences([in_text])[0]
    sequence = pad_sequences([sequence], max_length)

    y_pred = model.predict([feature, sequence])
    y_pred = np.argmax(y_pred)

    word = idx_to_word(y_pred, tokenizer)

    if word is None:
        break

    in_text += " " + word

    if word == 'endseq':
        break

return in_text
```

```
[26]: samples = test.sample(15)
       samples.reset_index(drop=True, inplace=True)
```

```
[27]:  
for index, record in samples.iterrows():  
  
    img = load_img(os.path.join(image_path, record['image']), target_size=(224,224))  
    img = img_to_array(img)  
    img = img/255.  
  
    caption = predict_caption(caption_model, record['image'], tokenizer, max_length, features)  
    samples.loc[index, 'caption'] = caption
```

```

1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 19ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 19ms/step
1/1 [=====] - 0s 19ms/step
1/1 [=====] - 0s 19ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 19ms/step
1/1 [=====] - 0s 24ms/step
1/1 [=====] - 0s 22ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 19ms/step
1/1 [=====] - 0s 19ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 19ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 19ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 19ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 19ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 24ms/step
1/1 [=====] - 0s 22ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 23ms/step
1/1 [=====] - 0s 23ms/step
1/1 [=====] - 0s 22ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 23ms/step
1/1 [=====] - 0s 22ms/step
1/1 [=====] - 0s 22ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 22ms/step
1/1 [=====] - 0s 19ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 20ms/step
1/1 [=====] - 0s 20ms/step

```

[28]: `display_images(samples)`

startseq group of people are playing in the grass endseq



startseq two dogs are playing with the ball endseq



startseq two children are playing in the water endseq



startseq dog is running through the grass endseq



startseq dog is running in the grass endseq



startseq man in red shirt is standing on the street endseq



startseq little girl is playing in the playground endseq



startseq man in black jacket is standing on the street endseq



startseq boy in blue shirt is jumping into the air endseq



startseq group of people are standing in front of crowd endseq



startseq the man is jumping into the water endseq



startseq man is jumping over rock endseq



startseq boy in blue shirt is playing in the water endseq



startseq black dog runs through the grass endseq



startseq person in red jacket is riding bike on the hill endseq



**Conclusion:** This may not be the best performing model, but the objective of this kernel is to give a gist of how Image Captioning problems can be approached. In the future work of this kernel Attention model training and BLEU Score assessment will be performed.

+ Code

+ Markdown

