

```
1 import tensorflow as tf
2
3 # You'll generate plots of attention in order to see which parts of an image
4 # our model focuses on during captioning
5 import matplotlib.pyplot as plt
6
7 # Scikit-learn includes many helpful utilities
8 from sklearn.model_selection import train_test_split
9 from sklearn.utils import shuffle
10
11 import re
12 import numpy as np
13 import os
14 import time
15 import json
16 from glob import glob
17 from PIL import Image
18 import pickle

1 # Download caption annotation files
2 annotation_folder = '/annotations/'
3 if not os.path.exists(os.path.abspath('.') + annotation_folder):
4     annotation_zip = tf.keras.utils.get_file('captions.zip',
5                                              cache_subdir=os.path.abspath('.'),
6                                              origin = 'http://images.cocodataset.org/captions\_train2014.zip',
7                                              extract = True)
8     annotation_file = os.path.dirname(annotation_zip) + '/annotations/captions_train'
9     os.remove(annotation_zip)
10
11 # Download image files
12 image_folder = '/train2014/'
13 if not os.path.exists(os.path.abspath('.') + image_folder):
14     image_zip = tf.keras.utils.get_file('train2014.zip',
15                                         cache_subdir=os.path.abspath('.'),
16                                         origin = 'http://images.cocodataset.org/zips/train2014.zip',
17                                         extract = True)
18     PATH = os.path.dirname(image_zip) + image_folder
19     os.remove(image_zip)
20 else:
21     PATH = os.path.abspath('.') + image_folder
```

```
↳ Downloading data from http://images.cocodataset.org/annotations/annotations\_train2014.zip
252878848/252872794 [=====] - 5s 0us/step
Downloading data from http://images.cocodataset.org/zips/train2014.zip
13510574080/13510573713 [=====] - 326s 0us/step
```

```
1 # Read the json file
2 with open(annotation_file, 'r') as f:
3     annotations = json.load(f)
4
5 # Store captions and image names in vectors
6 all_captions = []
7 all_img_name_vector = []
```

```

8
9 for annot in annotations['annotations']:
10    caption = '<start> ' + annot['caption'] + ' <end>'
11    image_id = annot['image_id']
12    full_coco_image_path = PATH + 'COCO_train2014_' + '%012d.jpg' % (image_id)
13
14    all_img_name_vector.append(full_coco_image_path)
15    all_captions.append(caption)
16
17 # Shuffle captions and image_names together
18 # Set a random state
19 train_captions, img_name_vector = shuffle(all_captions,
20                                             all_img_name_vector,
21                                             random_state=1)
22
23 # Select the first 30000 captions from the shuffled set
24 num_examples = 30000
25 train_captions = train_captions[:num_examples]
26 img_name_vector = img_name_vector[:num_examples]

1 len(train_captions), len(all_captions)


```

↳ (30000, 414113)

```

1 def load_image(image_path):
2     img = tf.io.read_file(image_path)
3     img = tf.image.decode_jpeg(img, channels=3)
4     img = tf.image.resize(img, (299, 299))
5     img = tf.keras.applications.inception_v3.preprocess_input(img)
6     return img, image_path

1 image_model = tf.keras.applications.InceptionV3(include_top=False,
2                                                 weights='imagenet')
3 new_input = image_model.input
4 hidden_layer = image_model.layers[-1].output
5
6 image_features_extract_model = tf.keras.Model(new_input, hidden_layer)


```

↳ Downloading data from https://storage.googleapis.com/tensorflow/keras-applications/inception_v3_weights_tf_dim_ordering_tf_kernels.h5
87916544/87910968 [=====] - 1s 0us/step

1 !pip install tqdm

↳ Requirement already satisfied: tqdm in /usr/local/lib/python3.6/dist-packages

```

1 from tqdm import tqdm
2
3

1 # Get unique images
2 encode_train = sorted(set(img_name_vector))
3
4 # Feel free to change batch_size according to your system configuration


```

```

5 image_dataset = tf.data.Dataset.from_tensor_slices(encode_train)
6 image_dataset = image_dataset.map(
7     load_image, num_parallel_calls=tf.data.experimental.AUTOTUNE).batch(16)
8
9 for img, path in tqdm(image_dataset):
10    batch_features = image_features_extract_model(img)
11    batch_features = tf.reshape(batch_features,
12                                (batch_features.shape[0], -1, batch_features.shape
13
14    for bf, p in zip(batch_features, path):
15        path_of_feature = p.numpy().decode("utf-8")
16        np.save(path_of_feature, bf.numpy())

```

↳ 1622it [05:34, 4.85it/s]

```

1 # Find the maximum length of any caption in our dataset
2 def calc_max_length(tensor):
3     return max(len(t) for t in tensor)

```

```

1 # Choose the top 5000 words from the vocabulary
2 top_k = 5000
3 tokenizer = tf.keras.preprocessing.text.Tokenizer(num_words=top_k,
4                                                 oov_token="",
5                                                 filters='!"#$%&()*+.,-/:;=?@\\'')
6 tokenizer.fit_on_texts(train_captions)
7 train_seqs = tokenizer.texts_to_sequences(train_captions)

```

```

1 tokenizer.word_index['<pad>'] = 0
2 tokenizer.index_word[0] = '<pad>'

```

```

1 # Create the tokenized vectors
2 train_seqs = tokenizer.texts_to_sequences(train_captions)

```

```

1 # Pad each vector to the max_length of the captions
2 # If you do not provide a max_length value, pad_sequences calculates it automatically
3 cap_vector = tf.keras.preprocessing.sequence.pad_sequences(train_seqs, padding=''

```

```

1 # Calculates the max_length, which is used to store the attention weights
2 max_length = calc_max_length(train_seqs)

```

```

1 # Create training and validation sets using an 80-20 split
2 img_name_train, img_name_val, cap_train, cap_val = train_test_split(img_name_vec
3                                         cap_vector,
4                                         test_size=0.
5                                         random_state

```

```

1 len(img_name_train), len(cap_train), len(img_name_val), len(cap_val)

```

↳ (24000, 24000, 6000, 6000)

```

1 # Feel free to change these parameters according to your system's configuration

```

```
1
2
3 BATCH_SIZE = 64
4 BUFFER_SIZE = 1000
5 embedding_dim = 256
6 units = 512
7 vocab_size = top_k + 1
8 num_steps = len(img_name_train) // BATCH_SIZE
9 # Shape of the vector extracted from InceptionV3 is (64, 2048)
10 # These two variables represent that vector shape
11 features_shape = 2048
12 attention_features_shape = 64

1 # Load the numpy files
2 def map_func(img_name, cap):
3     img_tensor = np.load(img_name.decode('utf-8')+'.npy')
4     return img_tensor, cap

1 dataset = tf.data.Dataset.from_tensor_slices((img_name_train, cap_train))
2
3 # Use map to load the numpy files in parallel
4 dataset = dataset.map(lambda item1, item2: tf.numpy_function(
5     map_func, [item1, item2], [tf.float32, tf.int32]),
6     num_parallel_calls=tf.data.experimental.AUTOTUNE)
7
8 # Shuffle and batch
9 dataset = dataset.shuffle(BUFFER_SIZE).batch(BATCH_SIZE)
10 dataset = dataset.prefetch(buffer_size=tf.data.experimental.AUTOTUNE)

1 class BahdanauAttention(tf.keras.Model):
2     def __init__(self, units):
3         super(BahdanauAttention, self).__init__()
4         self.W1 = tf.keras.layers.Dense(units)
5         self.W2 = tf.keras.layers.Dense(units)
6         self.V = tf.keras.layers.Dense(1)
7
8     def call(self, features, hidden):
9         # features(CNN_encoder output) shape == (batch_size, 64, embedding_dim)
10
11        # hidden shape == (batch_size, hidden_size)
12        # hidden_with_time_axis shape == (batch_size, 1, hidden_size)
13        hidden_with_time_axis = tf.expand_dims(hidden, 1)
14
15        # score shape == (batch_size, 64, hidden_size)
16        score = tf.nn.tanh(self.W1(features) + self.W2(hidden_with_time_axis))
17
18        # attention_weights shape == (batch_size, 64, 1)
19        # you get 1 at the last axis because you are applying score to self.V
20        attention_weights = tf.nn.softmax(self.V(score), axis=1)
21
22        # context_vector shape after sum == (batch_size, hidden_size)
23        context_vector = attention_weights * features
24        context_vector = tf.reduce_sum(context_vector, axis=1)
25
26        return context_vector, attention_weights
```

```
1 class CNN_Encoder(tf.keras.Model):
2     # Since you have already extracted the features and dumped it using pickle
3     # This encoder passes those features through a Fully connected layer
4     def __init__(self, embedding_dim):
5         super(CNN_Encoder, self).__init__()
6         # shape after fc == (batch_size, 64, embedding_dim)
7         self.fc = tf.keras.layers.Dense(embedding_dim)
8
9     def call(self, x):
10        x = self.fc(x)
11        x = tf.nn.relu(x)
12        return x
13
14
1 class RNN_Decoder(tf.keras.Model):
15     def __init__(self, embedding_dim, units, vocab_size):
16         super(RNN_Decoder, self).__init__()
17         self.units = units
18
19         self.embedding = tf.keras.layers.Embedding(vocab_size, embedding_dim)
20         self.gru = tf.keras.layers.GRU(self.units,
21                                       return_sequences=True,
22                                       return_state=True,
23                                       recurrent_initializer='glorot_uniform')
24         self.fc1 = tf.keras.layers.Dense(self.units)
25         self.fc2 = tf.keras.layers.Dense(vocab_size)
26
27         self.attention = BahdanauAttention(self.units)
28
29     def call(self, x, features, hidden):
30         # defining attention as a separate model
31         context_vector, attention_weights = self.attention(features, hidden)
32
33         # x shape after passing through embedding == (batch_size, 1, embedding_dim)
34         x = self.embedding(x)
35
36         # x shape after concatenation == (batch_size, 1, embedding_dim + hidden_size)
37         x = tf.concat([tf.expand_dims(context_vector, 1), x], axis=-1)
38
39         # passing the concatenated vector to the GRU
40         output, state = self.gru(x)
41
42         # shape == (batch_size, max_length, hidden_size)
43         x = self.fc1(output)
44
45         # x shape == (batch_size * max_length, hidden_size)
46         x = tf.reshape(x, (-1, x.shape[2]))
47
48         # output shape == (batch_size * max_length, vocab)
49         x = self.fc2(x)
50
51         return x, state, attention_weights
52
53     def reset_state(self, batch_size):
```

```

41     return tf.zeros((batch_size, self.units))

1 encoder = CNN_Encoder(embedding_dim)
2 decoder = RNN_Decoder(embedding_dim, units, vocab_size)

1 optimizer = tf.keras.optimizers.Adam()
2 loss_object = tf.keras.losses.SparseCategoricalCrossentropy(
3     from_logits=True, reduction='none')
4
5 def loss_function(real, pred):
6     mask = tf.math.logical_not(tf.math.equal(real, 0))
7     loss_ = loss_object(real, pred)
8
9     mask = tf.cast(mask, dtype=loss_.dtype)
10    loss_ *= mask
11
12    return tf.reduce_mean(loss_)

1 checkpoint_path = "./checkpoints/train"
2 ckpt = tf.train.Checkpoint(encoder=encoder,
3                             decoder=decoder,
4                             optimizer = optimizer)
5 ckpt_manager = tf.train.CheckpointManager(ckpt, checkpoint_path, max_to_keep=5)

1 start_epoch = 0
2 if ckpt_manager.latest_checkpoint:
3     start_epoch = int(ckpt_manager.latest_checkpoint.split('-')[-1])
4     # restoring the latest checkpoint in checkpoint_path
5     ckpt.restore(ckpt_manager.latest_checkpoint)

```

Training

```

1 # adding this in a separate cell because if you run the training cell
2 # many times, the loss_plot array will be reset
3 loss_plot = []

1 @tf.function
2 def train_step(img_tensor, target):
3     loss = 0
4
5     # initializing the hidden state for each batch
6     # because the captions are not related from image to image
7     hidden = decoder.reset_state(batch_size=target.shape[0])
8
9     dec_input = tf.expand_dims([tokenizer.word_index['<start>']] * target.shape[0]
10
11    with tf.GradientTape() as tape:
12        features = encoder(img_tensor)
13
14        for i in range(1, target.shape[1]):
15            # passing the features through the decoder
16            predictions, hidden, _ = decoder(dec_input, features, hidden)

```

```
17
18     loss += loss_function(target[:, i], predictions)
19
20     # using teacher forcing
21     dec_input = tf.expand_dims(target[:, i], 1)
22
23 total_loss = (loss / int(target.shape[1]))
24
25 trainable_variables = encoder.trainable_variables + decoder.trainable_variables
26
27 gradients = tape.gradient(loss, trainable_variables)
28
29 optimizer.apply_gradients(zip(gradients, trainable_variables))
30
31 return loss, total_loss

1 EPOCHS = 20
2
3 for epoch in range(start_epoch, EPOCHS):
4     start = time.time()
5     total_loss = 0
6
7     for (batch, (img_tensor, target)) in enumerate(dataset):
8         batch_loss, t_loss = train_step(img_tensor, target)
9         total_loss += t_loss
10
11     if batch % 100 == 0:
12         print ('Epoch {} Batch {} Loss {:.4f}'.format(
13             epoch + 1, batch, batch_loss.numpy() / int(target.shape[1])))
14     # storing the epoch end loss value to plot later
15     loss_plot.append(total_loss / num_steps)
16
17     if epoch % 5 == 0:
18         ckpt_manager.save()
19
20     print ('Epoch {} Loss {:.6f}'.format(epoch + 1,
21                                         total_loss/num_steps))
22     print ('Time taken for 1 epoch {} sec\n'.format(time.time() - start))
```



```
Epoch 1 Batch 0 Loss 1.9629
Epoch 1 Batch 100 Loss 1.0641
Epoch 1 Batch 200 Loss 1.0097
Epoch 1 Batch 300 Loss 0.9003
Epoch 1 Loss 1.039736
Time taken for 1 epoch 371.2840042114258 sec
```

```
Epoch 2 Batch 0 Loss 0.7678
Epoch 2 Batch 100 Loss 0.8646
Epoch 2 Batch 200 Loss 0.9211
Epoch 2 Batch 300 Loss 0.7804
Epoch 2 Loss 0.795022
Time taken for 1 epoch 352.27037048339844 sec
```

```
Epoch 3 Batch 0 Loss 0.7298
Epoch 3 Batch 100 Loss 0.7947
Epoch 3 Batch 200 Loss 0.6727
Epoch 3 Batch 300 Loss 0.7272
Epoch 3 Loss 0.723231
Time taken for 1 epoch 349.2909016609192 sec
```

```
Epoch 4 Batch 0 Loss 0.6679
Epoch 4 Batch 100 Loss 0.6928
Epoch 4 Batch 200 Loss 0.6878
Epoch 4 Batch 300 Loss 0.6621
Epoch 4 Loss 0.677160
Time taken for 1 epoch 348.0951991081238 sec
```

```
Epoch 5 Batch 0 Loss 0.6100
Epoch 5 Batch 100 Loss 0.6671
Epoch 5 Batch 200 Loss 0.5985
Epoch 5 Batch 300 Loss 0.5920
Epoch 5 Loss 0.640795
Time taken for 1 epoch 347.40416526794434 sec
```

```
Epoch 6 Batch 0 Loss 0.6440
Epoch 6 Batch 100 Loss 0.5976
Epoch 6 Batch 200 Loss 0.6401
Epoch 6 Batch 300 Loss 0.6115
Epoch 6 Loss 0.607960
Time taken for 1 epoch 346.4037923812866 sec
```

```
Epoch 7 Batch 0 Loss 0.5286
Epoch 7 Batch 100 Loss 0.6015
Epoch 7 Batch 200 Loss 0.6223
Epoch 7 Batch 300 Loss 0.5785
Epoch 7 Loss 0.577220
Time taken for 1 epoch 347.7673637866974 sec
```

```
Epoch 8 Batch 0 Loss 0.5193
Epoch 8 Batch 100 Loss 0.5708
Epoch 8 Batch 200 Loss 0.5374
Epoch 8 Batch 300 Loss 0.5351
Epoch 8 Loss 0.546481
Time taken for 1 epoch 346.2396740913391 sec
```

```
Epoch 9 Batch 0 Loss 0.5695
Epoch 9 Batch 100 Loss 0.5326
Epoch 9 Batch 200 Loss 0.5300
Epoch 9 Batch 300 Loss 0.5603
Epoch 9 Loss 0.516191
```

```
Time taken for 1 epoch 346.0829641819 sec
```

```
Epoch 10 Batch 0 Loss 0.4980  
Epoch 10 Batch 100 Loss 0.4978  
Epoch 10 Batch 200 Loss 0.4818  
Epoch 10 Batch 300 Loss 0.4839  
Epoch 10 Loss 0.487158
```

```
Time taken for 1 epoch 346.4885756969452 sec
```

```
Epoch 11 Batch 0 Loss 0.4478  
Epoch 11 Batch 100 Loss 0.5002  
Epoch 11 Batch 200 Loss 0.4418  
Epoch 11 Batch 300 Loss 0.4151  
Epoch 11 Loss 0.458268
```

```
Time taken for 1 epoch 346.2432270050049 sec
```

```
Epoch 12 Batch 0 Loss 0.4603  
Epoch 12 Batch 100 Loss 0.4438  
Epoch 12 Batch 200 Loss 0.3858  
Epoch 12 Batch 300 Loss 0.4008  
Epoch 12 Loss 0.429968
```

```
Time taken for 1 epoch 345.3726074695587 sec
```

```
Epoch 13 Batch 0 Loss 0.4132  
Epoch 13 Batch 100 Loss 0.4539  
Epoch 13 Batch 200 Loss 0.3743  
Epoch 13 Batch 300 Loss 0.3916  
Epoch 13 Loss 0.403569
```

```
Time taken for 1 epoch 345.8834364414215 sec
```

```
Epoch 14 Batch 0 Loss 0.4062  
Epoch 14 Batch 100 Loss 0.3725  
Epoch 14 Batch 200 Loss 0.3627  
Epoch 14 Batch 300 Loss 0.3957  
Epoch 14 Loss 0.377096
```

```
Time taken for 1 epoch 345.52081847190857 sec
```

```
Epoch 15 Batch 0 Loss 0.3490  
Epoch 15 Batch 100 Loss 0.3553  
Epoch 15 Batch 200 Loss 0.3811  
Epoch 15 Batch 300 Loss 0.3235  
Epoch 15 Loss 0.353224
```

```
Time taken for 1 epoch 345.56896018981934 sec
```

```
Epoch 16 Batch 0 Loss 0.4001  
Epoch 16 Batch 100 Loss 0.3422  
Epoch 16 Batch 200 Loss 0.3360  
Epoch 16 Batch 300 Loss 0.3295  
Epoch 16 Loss 0.330572
```

```
Time taken for 1 epoch 345.6349630355835 sec
```

```
Epoch 17 Batch 0 Loss 0.3651  
Epoch 17 Batch 100 Loss 0.3014  
Epoch 17 Batch 200 Loss 0.3094  
Epoch 17 Batch 300 Loss 0.3254  
Epoch 17 Loss 0.307942
```

```
Time taken for 1 epoch 346.03713154792786 sec
```

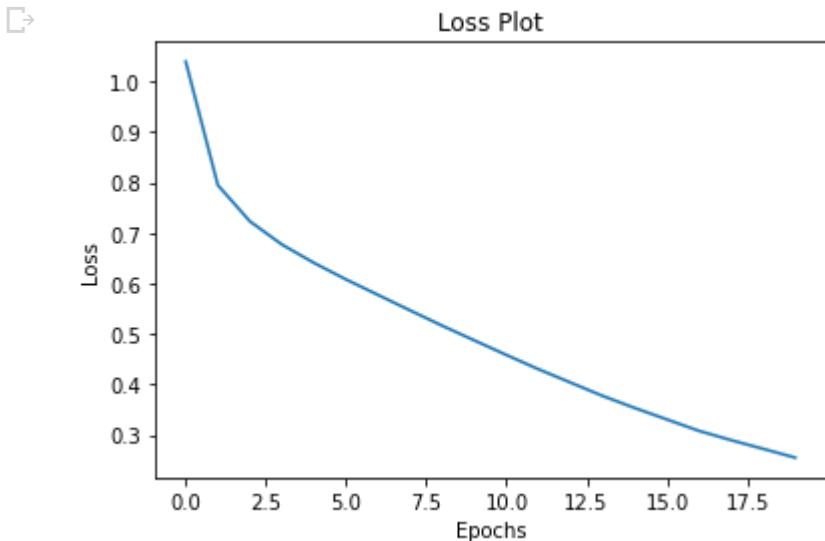
```
Epoch 18 Batch 0 Loss 0.3328  
Epoch 18 Batch 100 Loss 0.2998  
Epoch 18 Batch 200 Loss 0.2927  
Epoch 18 Batch 300 Loss 0.2962
```

```
Epoch 18 Loss 0.289562
Time taken for 1 epoch 345.2148702144623 sec
```

```
Epoch 19 Batch 0 Loss 0.2967
Epoch 19 Batch 100 Loss 0.2594
Epoch 19 Batch 200 Loss 0.2840
Epoch 19 Batch 300 Loss 0.2857
Epoch 19 Loss 0.272336
Time taken for 1 epoch 345.64116311073303 sec
```

```
Epoch 20 Batch 0 Loss 0.3181
Epoch 20 Batch 100 Loss 0.2437
Epoch 20 Batch 200 Loss 0.2525
Epoch 20 Batch 300 Loss 0.2457
Epoch 20 Loss 0.254878
Time taken for 1 epoch 344.66556310653687 sec
```

```
1 plt.plot(loss_plot)
2 plt.xlabel('Epochs')
3 plt.ylabel('Loss')
4 plt.title('Loss Plot')
5 plt.show()
```



Caption

```
1 def evaluate(image):
2     attention_plot = np.zeros((max_length, attention_features_shape))
3
4     hidden = decoder.reset_state(batch_size=1)
5
6     temp_input = tf.expand_dims(load_image(image)[0], 0)
7     img_tensor_val = image_features_extract_model(temp_input)
8     img_tensor_val = tf.reshape(img_tensor_val, (img_tensor_val.shape[0], -1, im
9
10    features = encoder(img_tensor_val)
11
12    dec_input = tf.expand_dims([tokenizer.word_index['<start>']], 0)
13    result = [1
```

```

14
15     for i in range(max_length):
16         predictions, hidden, attention_weights = decoder(dec_input, features, hi
17
18         attention_plot[i] = tf.reshape(attention_weights, (-1, )).numpy()
19
20         predicted_id = tf.random.categorical(predictions, 1)[0][0].numpy()
21         result.append(tokenizer.index_word[predicted_id])
22
23         if tokenizer.index_word[predicted_id] == '<end>':
24             return result, attention_plot
25
26         dec_input = tf.expand_dims([predicted_id], 0)
27
28     attention_plot = attention_plot[:len(result), :]
29     return result, attention_plot

1 def plot_attention(image, result, attention_plot):
2     temp_image = np.array(Image.open(image))
3
4     fig = plt.figure(figsize=(10, 10))
5
6     len_result = len(result)
7     for l in range(len_result):
8         temp_att = np.resize(attention_plot[l], (8, 8))
9         ax = fig.add_subplot(len_result//2, len_result//2, l+1)
10        ax.set_title(result[l])
11        img = ax.imshow(temp_image)
12        ax.imshow(temp_att, cmap='gray', alpha=0.6, extent=img.get_extent())
13
14    plt.tight_layout()
15    plt.show()

1 # captions on the validation set
2 rid = np.random.randint(0, len(img_name_val))
3 image = img_name_val[rid]
4 real_caption = ' '.join([tokenizer.index_word[i] for i in cap_val[rid] if i not
5 result, attention_plot = evaluate(image)
6
7 print ('Real Caption:', real_caption)
8 print ('Prediction Caption:', ' '.join(result))
9 #plot_attention(image, result, attention_plot)
10

⇨ Real Caption: <start> a tennis player waiting for the ball to be served <end>
Prediction Caption: a player dressed in red and <unk> a tennis <end>
```

Try

```

1 image_url = 'https://upload.wikimedia.org/wikipedia/commons/thumb/4/4b/Domestic_
2 image_extension = image_url[-4:]
3 image_path = tf.keras.utils.get_file('image15'+image_extension,
4                                     origin=image_url)
```

```
5  
6 result, attention_plot = evaluate(image_path)  
7 print ('Prediction Caption:', ' '.join(result))  
8 plot_attention(image_path, result, attention_plot)  
9 # opening the image  
10 Image.open(image_path)
```

□→

Downloading data from [https://upload.wikimedia.org/wikipedia/commons/thumb/4/41/475136/474072\[=====\]](https://upload.wikimedia.org/wikipedia/commons/thumb/4/41/475136/474072[=====]) [=====] – 0s 1us/step
Prediction Caption: a close up image of a cat cat by a pen <end>

