

# Task 3 Report: Image Captioning with AI

## Introduction

For Task 3 of my CodSoft internship, I used an AI model to generate a caption for an image I uploaded. The image features a blue BMW car parked next to a grassy area with autumn trees.

## Image Description

The image shows a blue BMW car, likely a 3 Series model, parked on a paved surface beside a grassy area. The car has a sleek design with a black kidney grille, sharp LED headlights, and the BMW logo in the center. The background includes green grass and trees with autumn-colored leaves, suggesting a fall setting.

## AI-Generated Caption

Using a pre-trained image captioning model in Google Colab, I generated this caption for my image:

**"a blue car is parked on the street."**

## How I Did It

I used Google Colab to run a notebook with a pre-trained image captioning model. I uploaded my photo of the BMW car, ran the code, and the model produced a caption. The code and output are shown below.

## Screenshots

- **Figure 1: Screenshot of the code (Part 1).**

This shows the first part of the code I ran.

A screenshot of the Google Colab code editor interface. The code is written in Python and is organized into six numbered steps. Step 1: Install the necessary tools (only need to run this once). Step 2: Load the tools we need. Step 3: Set up the pre-trained image captioning model. Step 4: Define a function to generate captions. Step 5: Upload your image and generate a caption. Step 6: Get the caption for the uploaded image. The code includes imports for transformers, PIL, and requests, and uses VisionEncoderDecoderModel, ViTFeatureExtractor, and AutoTokenizer. It also includes a function to open and process an image file. The interface shows a sidebar with icons for file explorer, code editor, and terminal. The top bar shows RAM and Disk usage, and the bottom bar shows the time as 10:43 PM and the Python version as Python 3.

```
# Step 1: Install the necessary tools (only need to run this once)
!pip install transformers
!pip install pillow

# Step 2: Load the tools we need
from transformers import VisionEncoderDecoderModel, ViTFeatureExtractor, AutoTokenizer
from PIL import Image
import requests

# Step 3: Set up the pre-trained image captioning model
model = VisionEncoderDecoderModel.from_pretrained("nlpconnect/vit-gpt2-image-captioning")
feature_extractor = ViTFeatureExtractor.from_pretrained("nlpconnect/vit-gpt2-image-captioning")
tokenizer = AutoTokenizer.from_pretrained("nlpconnect/vit-gpt2-image-captioning")

# Step 4: Define a function to generate captions
def generate_caption(image_path):
    # Open the image from your file
    image = Image.open(image_path).convert("RGB")

    # Prepare the image for the model
    pixel_values = feature_extractor(images=image, return_tensors="pt").pixel_values

    # Generate the caption
    output_ids = model.generate(pixel_values)
    caption = tokenizer.decode(output_ids[0], skip_special_tokens=True)

    return caption

# Step 5: Upload your image and generate a caption
from google.colab import files
uploaded = files.upload() # This opens a window to upload your image

# Step 6: Get the caption for the uploaded image
```



- **Figure 4: Screenshot of the output (Part 2).**

This continues the output.

```

    },
    "attention_probs_dropout_prob": 0.0,
    "encoder_stride": 16,
    "hidden_act": "gelu",
    "hidden_dropout_prob": 0.0,
    "hidden_size": 768,
    "image_size": 224,
    "initializer_range": 0.02,
    "intermediate_size": 3072,
    "layer_norm_eps": 1e-12,
    "model_type": "vit",
    "num_attention_heads": 12,
    "num_channels": 3,
    "num_hidden_layers": 12,
    "patch_size": 16,
    "pooler_act": "tanh",
    "pooler_output_size": 768,
    "qkv_bias": true,
    "torch_dtype": "float32",
    "transformers_version": "4.51.3"
  }

  Config of the decoder: <class 'transformers.models.gpt2.modeling_gpt2.GPT2LMHeadModel'> is overwritten by shared decoder config: GPT2Config {
    "activation_function": "gelu_new",
    "add_cross_attention": true,
    "architectures": [
      "GPT2LMHeadModel"
    ],
    "attn_pdrop": 0.1,
    "bos_token_id": 50256,
    "decoder_start_token_id": 50256,
    "embd_pdrop": 0.1,
    "eos_token_id": 50256,
    "initializer_range": 0.02,
  }

```

- **Figure 5: Screenshot of the output (Part 3).**

This shows more of the output.

```

    "is_decoder": true,
    "layer_norm_epsilon": 1e-05,
    "model_type": "gpt2",
    "n_ctx": 1024,
    "n_embd": 768,
    "n_head": 12,
    "n_inner": null,
    "n_layer": 12,
    "n_positions": 1024,
    "pad_token_id": 50256,
    "reorder_and_upcast_attn": false,
    "resid_pdrop": 0.1,
    "scale_attn_by_inverse_layer_idx": false,
    "scale_attn_weights": true,
    "summary_activation": null,
    "summary_first_dropout": 0.1,
    "summary_proj_to_labels": true,
    "summary_type": "cls_index",
    "summary_use_proj": true,
    "task_specific_params": {
      "text-generation": {
        "do_sample": true,
        "max_length": 50
      }
    },
    "torch_dtype": "float32",
    "transformers_version": "4.51.3",
    "use_cache": true,
    "vocab_size": 50257
  }

  preprocessor_config.json: 100% ██████████ 228/228 [00:00<00:00, 15.0KB/s]
  /usr/local/lib/python3.11/dist-packages/transformers/models/vit/feature_extraction_vit.py:28: FutureWarning: The class ViTFeatureExtractor is deprecated and will be removed in versio
  warnings.warn(
  tokenizer_config.json: 100% ██████████ 241/241 [00:00<00:00, 10.7KB/s]

```

- **Figure 6: Screenshot of the output (Part 4).**

This shows the final output with the caption.



```
preprocessor_config.json: 100% 228/228 [00:00<00:00, 15.0kB/s]
/usr/local/lib/python3.11/dist-packages/transformers/models/vit/feature_extraction_vit.py:28: FutureWarning: The class ViTFeatureExtractor is deprecated and will be removed in version 4.42.0. Please use ViTImageProcessor instead.
  warnings.warn(
tokenizer_config.json: 100% 241/241 [00:00<00:00, 19.7kB/s]
vocab.json: 100% 798k/798k [00:00<00:00, 2.90MB/s]
merges.txt: 100% 456k/456k [00:00<00:00, 7.45MB/s]
tokenizer.json: 100% 1.36M/1.36M [00:00<00:00, 10.8MB/s]
special_tokens_map.json: 100% 120/120 [00:00<00:00, 8.38kB/s]
Choose Files pexels-mike...-170811.jpg
• pexels-mikebirdy-170811.jpg(image/jpeg) - 1475919 bytes, last modified: 5/22/2025 - 100% done
Saving pexels-mikebirdy-170811.jpg to pexels-mikebirdy-170811.jpg
The attention mask is not set and cannot be inferred from input because pad token is same as eos token. As a consequence, you may observe unexpected behavior. Please pass your input's
attention_mask explicitly to the model. We strongly recommend passing in an 'attention_mask' since your input_ids may be padded. See https://huggingface.co/docs/transformers/troubleshooting#incorrect-output-when-padding-tokens-are-used
You may ignore this warning if your 'pad_token_id' (50256) is identical to the 'bos_token_id' (50256), 'eos_token_id' (50256), or the 'sep_token_id' (None), and your input is not padded.
Image: pexels-mikebirdy-170811.jpg
Caption: a blue car is parked on the street

[ ] Start coding or generate with AI.
```

## Reflection

This task taught me how AI can create captions from images. My caption was accurate but missed some details like the trees. I enjoyed using Colab and seeing the model work with my photo.