

Setup

This first cell installs all the necessary libraries for this notebook. We are using `langchain`, `gradio`, `transformers`, `bs4`, `requests`, and `torch`.

```
[1] 1 # installing required libraries in my_env  
2 !pip install langchain==0.1.11 gradio==5.23.2 transformers==4.38.2 bs4==0.0.2 requests==2.31.0 torch==2.2.1
```

1- Image Captioning Model

Here we are importing the necessary libraries for the image captioning part of the notebook: `requests` to download images from the web, `PIL` (Pillow) to work with images, and `transformers` from Hugging Face to use a pre-trained image captioning model.

```
[1] ✓ 47s 1 import requests  
2 from PIL import Image  
3 from transformers import AutoProcessor, BlipForConditionalGeneration  
4  
5 # Load the pretrained processor and model  
6 processor = AutoProcessor.from_pretrained("Salesforce/blip-image-captioning-base")  
7 model = BlipForConditionalGeneration.from_pretrained("Salesforce/blip-image-captioning-base")
```

/usr/local/lib/python3.12/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning: The secret `HF_TOKEN` does not exist in your Colab secrets.

To authenticate with the Hugging Face Hub, create a token in your settings tab (<https://huggingface.co/settings/tokens>), set it as secret in your Google Colab and restart your notebook. You will be able to reuse this secret in all of your notebooks.

Please note that authentication is recommended but still optional to access public models or datasets.

warnings.warn(
preprocessor_config.json: 100% 287/287 [00:00<00:00, 20.7kB/s]

Using a slow image processor as `use_fast` is unset and a slow processor was saved with this model. `use_fast=True` will be the default behavior in v4.52, even if the model was trained with a fast processor.
tokenizer_config.json: 100% 506/506 [00:00<00:00, 30.5kB/s]

vocab.txt: 232k? [00:00<00:00, 10.3MB/s]

tokenizer.json: 711k? [00:00<00:00, 40.3MB/s]

special_tokens_map.json: 100% 125/125 [00:00<00:00, 12.9kB/s]

config.json: 4.56k? [00:00<00:00, 437kB/s]

pytorch_model.bin: 100% 990M/990M [00:07<00:00, 85.9MB/s]

This cell loads an image from a specified path. The code then converts the image to the RGB format, which is required by the image captioning model.

```
[2] ✓ 0s 1 img_path = "/content/sun.png"  
2  
3 # convert it into an RGB format  
4 image = Image.open(img_path).convert('RGB')  
5  
/usr/local/lib/python3.12/dist-packages/PIL/Image.py:1047: UserWarning: Palette images with Transparency expressed in bytes should be converted to RGBA images  
warnings.warn(
```

Here, we prepare the image and an optional text input for the model. For image captioning, a generic text like "the image about what" is used as a prompt. The `processor` converts the image and text into the format that the model expects.

```
[26] ✓ 0s 1 # You do not need a question for image captioning  
2 text = "This image shows"  
3 inputs = processor(images=image, text=text, return_tensors="pt")
```

The `generate` method of the model produces a sequence of tokens representing the caption. We limit the caption length to 50 tokens using `max_length=50`.

```
[27] 1 # Generate a caption for the image  
2 outputs = model.generate(**inputs, max_length=50)
```

After generating the caption tokens, this cell decodes them back into a human-readable text string. The `processor.decode` method is used for this purpose, and `skip_special_tokens=True` removes any special tokens added during the generation process.

```
[28] ✓ 0s 1 # Decode the generated tokens to text  
2 caption = processor.decode(outputs[0], skip_special_tokens=True)  
3 # Print the caption  
4 print(caption)
```

this image shows the sun with a smiley face

2- Gradio Interface

Import the necessary libraries for building the Gradio interface for our image captioning model.

```
[7] ✓ 9s 1 import gradio as gr  
2 import numpy as np  
3 from PIL import Image  
4 from transformers import AutoProcessor, BlipForConditionalGeneration
```

This function `caption_image` is the core logic for the Gradio interface. It takes a NumPy array representing an image as input. It converts the image to the correct format, uses the processor to prepare the input for the model, generates the caption using the model, and finally decodes the output into a readable string. This string is then returned as the output of the function.

```
[8] ✓ 0s 1 def caption_image(input_image: np.ndarray):  
2     # Convert numpy array to PIL Image and convert to RGB  
3     raw_image = Image.fromarray(input_image).convert('RGB')  
4  
5     # Process the image  
6     inputs = processor(raw_image, return_tensors="pt")  
7  
8     # Generate a caption for the image  
9     outputs = model.generate(**inputs)  
10  
11    # Decode the generated tokens to text and store it into `caption`  
12    caption = processor.decode(outputs[0], skip_special_tokens=True)  
13  
14    return caption
```

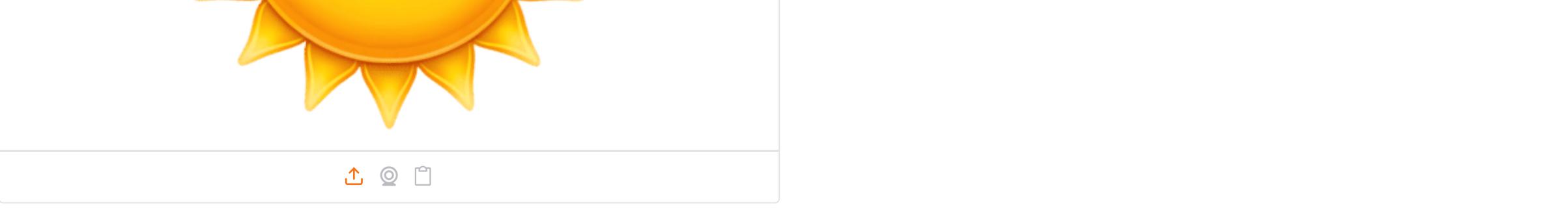
This cell defines the Gradio interface for the image captioning task. It uses the `caption_image` function as the core logic. The interface has an image input (`gr.Image()`) and a text output (`"text"`). We also provide a title and description for the web app.

```
[9] ✓ 0s 1 iface = gr.Interface(  
2     fn=caption_image,  
3     inputs=gr.Image(),  
4     outputs="text",  
5     title="Image Captioning",  
6     description="This is a simple web app for generating captions for images using a trained model."  
7 )
```

Finally, this cell launches the Gradio interface. Once executed, it will provide a public URL that you can use to access the image captioning web app.

```
[10] ✓ 1s 1 iface.launch()  
... It looks like you are running Gradio on a hosted Jupyter notebook, which requires `share=True`. Automatically setting `share=True` (you can turn this off by setting `share=False` in the launch() call)  
* Running on public URL: https://dfa1efea24afb23025.gradio.live
```

This share link expires in 1 week. For free permanent hosting and GPU upgrades, run `gradio deploy` from the terminal in the working directory to deploy to Hugging Face Space.



Next steps: [Deploy to Cloud Run](#)