

VisionGPT-2 Image Captioning Model

- *almost* built from scratch.
- pretrained weights loading via HF & timm
- dataset preparation from scratch as well.

Imports

```
1 ! pip install timm

→ Collecting timm
  Downloading timm-1.0.9-py3-none-any.whl.metadata (42 kB)
  ━━━━━━━━━━━━━━━━ 42.4/42.4 kB 3.3 MB/s eta 0:00:00
Requirement already satisfied: torch in /usr/local/lib/python3.10/dist-packages (from timm) (2.4.1+cu121)
Requirement already satisfied: torchvision in /usr/local/lib/python3.10/dist-packages (from timm) (0.19.1+cu121)
Requirement already satisfied: pyyaml in /usr/local/lib/python3.10/dist-packages (from timm) (6.0.2)
Requirement already satisfied: huggingface_hub in /usr/local/lib/python3.10/dist-packages (from timm) (0.24.7)
Requirement already satisfied: safetensors in /usr/local/lib/python3.10/dist-packages (from timm) (0.4.5)
Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from huggingface_hub->timm) (3.16.1)
Requirement already satisfied: fsspec>=2023.5.0 in /usr/local/lib/python3.10/dist-packages (from huggingface_hub->timm) (2024.6.1)
Requirement already satisfied: packaging>=20.9 in /usr/local/lib/python3.10/dist-packages (from huggingface_hub->timm) (24.1)
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from huggingface_hub->timm) (2.32.3)
Requirement already satisfied: tqdm>=4.42.1 in /usr/local/lib/python3.10/dist-packages (from huggingface_hub->timm) (4.66.5)
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.10/dist-packages (from huggingface_hub->timm) (<
Requirement already satisfied: sympy in /usr/local/lib/python3.10/dist-packages (from torch->timm) (1.13.3)
Requirement already satisfied: networkx in /usr/local/lib/python3.10/dist-packages (from torch->timm) (3.3)
Requirement already satisfied: jinja2 in /usr/local/lib/python3.10/dist-packages (from torch->timm) (3.1.4)
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (from torchvision->timm) (1.26.4)
Requirement already satisfied: pillow!=8.3.*,>=5.3.0 in /usr/local/lib/python3.10/dist-packages (from torchvision->timm) (10.4.0)
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/dist-packages (from jinja2->torch->timm) (2.1.5)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->huggingface_hub->timm)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->huggingface_hub->timm) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests->huggingface_hub->timm)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests->huggingface_hub->timm)
Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.10/dist-packages (from sympy->torch->timm) (1.3.0)
Downloading timm-1.0.9-py3-none-any.whl (2.3 MB)
  ━━━━━━━━━━━━━━━━ 2.3/2.3 MB 58.0 MB/s eta 0:00:00

Installing collected packages: timm
Successfully installed timm-1.0.9
```

```
1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import numpy as np
5 import pandas as pd
6 import matplotlib.pyplot as plt
7 from timm import create_model, list_models
8 from types import SimpleNamespace
9 from transformers import GPT2LMHeadModel, GPT2TokenizerFast, get_linear_schedule_with_warmup
10 import albumentations as A
11 from albumentations.pytorch import ToTensorV2
12 from PIL import Image
13 from pathlib import Path
14 from sklearn.model_selection import train_test_split
15 from torch.cuda.amp import GradScaler, autocast
16 from tqdm.auto import tqdm
17 import gc
18 import json

1 %env TOKENIZERS_PARALLELISM = false

→ env: TOKENIZERS_PARALLELISM=false
```

Dataset

- the dataset we're using is COCO 2017, it has about 500k samples, we will only use 150k samples.
- augmentations: PIL + albumentations. Fun fact: albumentations is a LOT faster than torchvision
- mean and std for ViT models is [0.5,0.5,0.5] unlike the standard ImageNet mean and std.

```

1 sample_tfms = [
2     A.HorizontalFlip(),
3     A.RandomBrightnessContrast(),
4     A.ColorJitter(),
5     A.ShiftScaleRotate(shift_limit=0.1, scale_limit=0.3, rotate_limit=45, p=0.5),
6     A.HueSaturationValue(p=0.3),
7 ]
8 train_tfms = A.Compose([
9     *sample_tfms,
10    A.Resize(224,224),
11    A.Normalize(mean=[0.5,0.5,0.5],std=[0.5,0.5,0.5],always_apply=True),
12    ToTensorV2()
13 ])
14 valid_tfms = A.Compose([
15    A.Resize(224,224),
16    A.Normalize(mean=[0.5,0.5,0.5],std=[0.5,0.5,0.5],always_apply=True),
17    ToTensorV2()
18 ])

1 tokenizer = GPT2TokenizerFast.from_pretrained('gpt2')
2 tokenizer.pad_token = tokenizer.eos_token
3 tokenizer.pad_token

✉ /usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_token.py:89: 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 :
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(
        tokenizer_config.json: 100%                                         26.0/26.0 [00:00<00:00, 1.51kB/s]
        vocab.json: 100%                                         1.04M/1.04M [00:00<00:00, 2.41MB/s]
        merges.txt: 100%                                         456k/456k [00:00<00:00, 712kB/s]
        tokenizer.json: 100%                                         1.36M/1.36M [00:00<00:00, 1.57MB/s]
        config.json: 100%                                         665/665 [00:00<00:00, 31.9kB/s]
/usr/local/lib/python3.10/dist-packages/transformers/tokenization_utils_base.py:1601: FutureWarning: `clean_up_tokenization_spaces`
    warnings.warn(
        '<|endoftext|>'

◀ ➤
```

1 tokenizer.encode_plus('bonjour! WISD caption ')

✉ {'input_ids': [4189, 73, 454, 0, 370, 1797, 35, 8305, 220], 'attention_mask': [1, 1, 1, 1, 1, 1, 1, 1, 1]}

```

1 class Dataset:
2     def __init__(self, df, tfms):
3         self.df = df
4         self.tfms = tfms
5     def __len__(self):
6         return len(self.df)
7     def __getitem__(self, idx):
8         sample = self.df.iloc[idx,:]
9         image = sample['image']
10        caption = sample['caption']
11        image = Image.open(image).convert('RGB')
12        image = np.array(image)
13        aug = self.tfms(image=image)
14        image = aug['image']
15        caption = f'{caption}<|endoftext|>'
16        input_ids = tokenizer(
17            caption,
18            truncation=True)['input_ids']
19        labels = input_ids.copy()
20        labels[:-1] = input_ids[1:]
21        return image, input_ids, labels

1 import os
2
3
4 from google.colab import files
5 files.upload() # Téléchargez 'kaggle.json' ici
6
7 !mkdir -p ~/.kaggle
8 !mv kaggle.json ~/.kaggle/
9 !chmod 600 ~/.kaggle/kaggle.json
10
11 # Installer l'outil Kaggle dans Colab
12 !pip install -q kaggle
```

13
14

Aucun fichier choisi Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving kaggle.json to kaggle (1).json

```
1 # Remplacer 'dataset' par l'URL ou l'ID du dataset Kaggle que vous voulez télécharger
2 # Par exemple, pour le dataset COCO 2017:
3 !kaggle datasets download -d awsaf49/coco-2017-dataset
4
5 # Extraire les fichiers du zip téléchargé
6 !unzip coco-2017-dataset.zip -d ./coco2017/
```

Le flux de sortie a été tronqué et ne contient que les 5000 dernières lignes.

```
inflating: ./coco2017/coco2017/val2017/000000000139.jpg
inflating: ./coco2017/coco2017/val2017/00000000285.jpg
inflating: ./coco2017/coco2017/val2017/00000000632.jpg
inflating: ./coco2017/coco2017/val2017/00000000724.jpg
inflating: ./coco2017/coco2017/val2017/00000000776.jpg
inflating: ./coco2017/coco2017/val2017/00000000785.jpg
inflating: ./coco2017/coco2017/val2017/00000000802.jpg
inflating: ./coco2017/coco2017/val2017/00000000872.jpg
inflating: ./coco2017/coco2017/val2017/00000000885.jpg
inflating: ./coco2017/coco2017/val2017/000000001000.jpg
inflating: ./coco2017/coco2017/val2017/000000001268.jpg
inflating: ./coco2017/coco2017/val2017/000000001296.jpg
inflating: ./coco2017/coco2017/val2017/000000001353.jpg
inflating: ./coco2017/coco2017/val2017/000000001425.jpg
inflating: ./coco2017/coco2017/val2017/000000001490.jpg
inflating: ./coco2017/coco2017/val2017/000000001503.jpg
inflating: ./coco2017/coco2017/val2017/000000001532.jpg
inflating: ./coco2017/coco2017/val2017/000000001584.jpg
inflating: ./coco2017/coco2017/val2017/000000001675.jpg
inflating: ./coco2017/coco2017/val2017/000000001761.jpg
inflating: ./coco2017/coco2017/val2017/000000001818.jpg
inflating: ./coco2017/coco2017/val2017/000000001993.jpg
inflating: ./coco2017/coco2017/val2017/000000002006.jpg
inflating: ./coco2017/coco2017/val2017/000000002149.jpg
inflating: ./coco2017/coco2017/val2017/000000002153.jpg
inflating: ./coco2017/coco2017/val2017/000000002157.jpg
inflating: ./coco2017/coco2017/val2017/000000002261.jpg
inflating: ./coco2017/coco2017/val2017/000000002299.jpg
inflating: ./coco2017/coco2017/val2017/000000002431.jpg
inflating: ./coco2017/coco2017/val2017/000000002473.jpg
inflating: ./coco2017/coco2017/val2017/000000002532.jpg
inflating: ./coco2017/coco2017/val2017/000000002587.jpg
inflating: ./coco2017/coco2017/val2017/000000002592.jpg
inflating: ./coco2017/coco2017/val2017/000000002685.jpg
inflating: ./coco2017/coco2017/val2017/000000002923.jpg
inflating: ./coco2017/coco2017/val2017/000000003156.jpg
inflating: ./coco2017/coco2017/val2017/000000003255.jpg
inflating: ./coco2017/coco2017/val2017/000000003501.jpg
inflating: ./coco2017/coco2017/val2017/000000003553.jpg
inflating: ./coco2017/coco2017/val2017/000000003661.jpg
inflating: ./coco2017/coco2017/val2017/000000003845.jpg
inflating: ./coco2017/coco2017/val2017/000000003934.jpg
inflating: ./coco2017/coco2017/val2017/000000004134.jpg
inflating: ./coco2017/coco2017/val2017/000000004395.jpg
inflating: ./coco2017/coco2017/val2017/000000004495.jpg
inflating: ./coco2017/coco2017/val2017/000000004765.jpg
inflating: ./coco2017/coco2017/val2017/000000004795.jpg
inflating: ./coco2017/coco2017/val2017/000000005001.jpg
inflating: ./coco2017/coco2017/val2017/000000005037.jpg
inflating: ./coco2017/coco2017/val2017/000000005060.jpg
inflating: ./coco2017/coco2017/val2017/000000005193.jpg
inflating: ./coco2017/coco2017/val2017/000000005477.jpg
inflating: ./coco2017/coco2017/val2017/000000005503.jpg
inflating: ./coco2017/coco2017/val2017/000000005529.jpg
inflating: ./coco2017/coco2017/val2017/000000005586.jpg
inflating: ./coco2017/coco2017/val2017/000000005600.jpg
inflating: ./coco2017/coco2017/val2017/000000005992.jpg
```

```
1
2
3 print(os.listdir('./coco2017'))
4
```

['coco2017']

```
1 # COCO 2017
2 base_path = Path('/content/coco2017/coco2017')
3 annot = base_path / 'annotations' / 'captions_train2017.json'
4 with open(annot, 'r') as f:
5     data = json.load(f)
```

```
6     data = data['annotations']
7
8 samples = []
9
10 for sample in data:
11     im = '%012d.jpg' % sample['image_id']
12     samples.append([im, sample['caption']])
13
14 df = pd.DataFrame(samples, columns=['image', 'caption'])
15 df['image'] = df['image'].apply(
16     lambda x: base_path / 'train2017' / x
17 )
18 df = df.sample(150_000)
19 df = df.reset_index(drop=True)
20 df.head()
```

	image	caption
0	/content/coco2017/coco2017/train2017/000000497...	A photo of a photo above the urinals in a publ...
1	/content/coco2017/coco2017/train2017/000000558...	There is a blue bird sitting on the person's h...
2	/content/coco2017/coco2017/train2017/000000205...	A photo of a person on a motorcycle on the road.
3	/content/coco2017/coco2017/train2017/000000439...	A man is walking with his horse along a grassy...
4	/content/coco2017/coco2017/train2017/000000171...	a close up of a tennis player swinging a racket

▼ some samples from the dataset

```
1 sampled_df = df.sample(n=20)
2 fig, axs = plt.subplots(10, 2, figsize=(20, 30))
3
4 for i, row in enumerate(sampled_df.iterrows()):
5     ax = axs[i // 2, i % 2]
6     image_path = row[1]['image']
7     caption = row[1]['caption']
8     image = Image.open(image_path)
9     ax.imshow(image)
10    ax.axis('off')
11    ax.set_title(caption)
12
13 plt.tight_layout()
14 plt.show()
```



A person on a snowboard riding over a hill.



Several skiers trace a path along a snowy slope.



The banana has many dark spots inside of it.



A giraffe laying in dirt area with a pavilion in the background.



A box of a variety of pastries on a table



A bathroom with toilet, bathtub, and towels in it.



A balloon of a giraffe next to a parking meter .



A colorful vase filled with a green plant.



A man carrying a surfboard near a beach with people on it.



A couple of guys playing with a frisbee



A red stop sign next to the side of a road.



a kid is riding a surfboard is riding on a wave



A bench beside a sidewalk in a grassy area.



A bathroom with a toilet, toilet tissue holding and a garbage can in it.



Three giraffes all interested in the activity at a zoo building.



A man is throwing a Frisbee in the air.



Memorial bench can be reached by foot or tram

a lantern sitting on top of a bar with bar stools in front of a sofa and coffee table filled with dishes

- ✓ you can also choose to use Flickr30k which has ~160k samples

```
1 # flickr30k
2 """
3 base_path = Path('/kaggle/input/flickr30k/flickr30k_images')
4 df = pd.read_csv('/kaggle/input/flickr30k/captions.txt', delimiter=',')
5 df.rename({'image_name': 'image', 'comment': 'caption'}, inplace=True, axis=1)
6 df['image'] = df['image'].map(lambda x: base_path / x.strip())
7 df['caption'] = df['caption'].map(lambda x: x.strip().lower())
8 df.head()
9 """

→ "nbase_path = Path('/kaggle/input/flickr30k/flickr30k_images')\ndf =\n    pd.read_csv('/kaggle/input/flickr30k/captions.txt', delimiter=',')\ndf.rename({'image_name': 'image', 'comment':\n        'caption'}, inplace=True, axis=1)\ndf['image'] = df['image'].map(lambda x: base_path / x.strip())\ndf['caption'] =\n    df['caption'].map(lambda x: x.strip().lower())\ndf.head()\n\n

1 train_df, val_df = train_test_split(df, test_size=0.1)
2 train_df.reset_index(drop=True, inplace=True)
3 val_df.reset_index(drop=True, inplace=True)
4 print(len(train_df), len(val_df))

→ 135000 15000

1 train_ds = Dataset(train_df, train_tfms)
2 val_ds = Dataset(val_df, valid_tfms)
```

✓ Custom collate function

- allows for dynamic padding so the model doesn't have to process `max_len` sequences which would be just filled with pad tokens
 - instead, every batch is padded according to the longest sequence in the batch

```
1 def collate_fn(batch):
2     image = [i[0] for i in batch]
3     input_ids = [i[1] for i in batch]
4     labels = [i[2] for i in batch]
5     image = torch.stack(image, dim=0)
6     input_ids = tokenizer.pad(
7         {'input_ids':input_ids},
8         padding='longest',
9         return_attention_mask=False,
10        return_tensors='pt'
11    )['input_ids']
12    labels = tokenizer.pad(
13        {'input_ids':labels},
14        padding='longest',
15        return_attention_mask=False,
16        return_tensors='pt'
17    )['input_ids']
18    mask = (input_ids!=tokenizer.pad_token_id).long()
19    labels[mask==0]=-100
20    return image, input_ids, labels
```

▼ How the data looks:

- every caption is a sequence of tokens, and as it is causal language modeling where the model predicts the next token, the labels are right shifted by 1 position.
 - every caption ends with the end of sentence token: eos_token (50256 : <|endoftext|>)
 - in GPT models, the pad tokens are same as the eos tokens, hence we also mask the pad tokens in the labels with -100 which are ignored by cross-entropy loss' default behaviour – check `collate_fn` to see how I masked them.

```
1 dl = torch.utils.data.DataLoader(train_ds,shuffle=True,batch_size=2,collate_fn=collate_fn)
2 _,c,l = next(iter(dl))
3 print(c[0])
4 print(l[0])

→ You're using a GPT2TokenizerFast tokenizer. Please note that with a fast tokenizer, using the `__call__` method is faster than using
  tensor([14945, 1180, 45613, 836, 5500, 1978, 287, 257, 3091, 13,
         50256])
  tensor([ 1180, 45613, 836, 5500, 1978, 287, 257, 3091, 13, 50256,
         -100])
```

▼ Model

-
- GPT2 on its own is a decoder-only model. So it understands only textual context. To create a model which also understands image context and applies it to the text, we use cross-attention.
 - Cross Attention: It is the second attention layer of the decoder block as per the [classic transformer encoder-decoder model](#)
 - The query in cross-attention is the output of the previous causal attention layer in the decoder block. The key and value are the outputs of the corresponding encoder block i.e. the image context.
 - Since GPT2 doesn't have a cross-attention layer, I coded it as `GPT2CrossAttention` and made it a part of the usual decoder `GPT2Block`.
 - The encoder I used was ViT-base with `patch_size=16` grabbed from `timm`.
 - ViT and GPT2 are extremely compatible with each other. They have the same amount of depth i.e. encoder, decoder blocks respectively. They have the same hidden size of 768 as well. So I didn't have to pool the encoder outputs to match the decoder hidden size. Less parameters yay!
 - GPT2 code I wrote was ofc based on what I learnt from [NanoGPT](#)

This is somewhat how the architecture looks now.

✓ Causal Attention Block

```

1 class GPT2Attention(nn.Module):
2     def __init__(self, config):
3         super().__init__()
4         self.embed_dim = config.embed_dim
5         self.n_heads = config.num_heads
6         assert self.embed_dim % self.n_heads == 0, 'embedding dimension must be divisible by number of heads'
7         self.head_size = self.embed_dim // self.n_heads
8         self.seq_len = config.seq_len
9
10        self.c_attn = nn.Linear(self.embed_dim, self.head_size * self.n_heads * 3, bias=True)
11        self.scale = self.head_size ** -0.5
12
13        self.register_buffer('mask', torch.tril(torch.ones(1, 1, self.seq_len, self.seq_len)))
14
15        self.c_proj = nn.Linear(self.embed_dim, self.embed_dim, bias=True)
16
17        self.attn_dropout = nn.Dropout(config.attention_dropout)
18        self.resid_dropout = nn.Dropout(config.residual_dropout)
19
20
21    def forward(self, x):
22        b, t, c = x.shape
23        # q, k, v shape individually: batch_size x seq_len x embed_dim
24        # we know that qk_t = q x k_t, where q=bxthead_dim, k_t=bxhead_dimxt
25        q, k, v = self.c_attn(x).chunk(3, dim=-1)
26        q = q.view(b, t, self.n_heads, self.head_size).permute(0, 2, 1, 3) # batch x n_heads x seq_len x head_dim
27        k = k.view(b, t, self.n_heads, self.head_size).permute(0, 2, 1, 3)
28        v = v.view(b, t, self.n_heads, self.head_size).permute(0, 2, 1, 3)
29
30        qk_t = (q @ k.transpose(-2, -1)) * self.scale
31        qk_t = qk_t.masked_fill(self.mask[:, :, :, :t] == 0, float('-inf'))
32        qk_t = F.softmax(qk_t, dim=-1)
33        weights = self.attn_dropout(qk_t)
34
35        attention = weights @ v # batch x n_heads x t x head_size
36        attention = attention.permute(0, 2, 1, 3).contiguous().view(b, t, c) # batch x t x embed_dim
37
38        out = self.c_proj(attention)
39        out = self.resid_dropout(out)
40
41        return out

```

✓ Cross Attention Block

- Notice how I initialized the parameters with mean=0. and std=0.02. This is an OpenAI trick they used while training GPT2 as well.

```

1 class GPT2CrossAttention(nn.Module):
2     def __init__(self, config):
3         super().__init__()
4         self.embed_dim = config.embed_dim
5         self.n_heads = config.num_heads
6         assert self.embed_dim % self.n_heads == 0, 'embedding dimension must be divisible by number of heads'
7         self.head_size = self.embed_dim // self.n_heads
8         self.seq_len = config.seq_len
9
10        self.q = nn.Linear(self.embed_dim, self.embed_dim)
11        self.k = nn.Linear(self.embed_dim, self.embed_dim)
12        self.v = nn.Linear(self.embed_dim, self.embed_dim)
13        self.scale = self.head_size ** -0.5
14
15        self.c_proj = nn.Linear(self.embed_dim, self.embed_dim, bias=True)
16
17        self.attn_dropout = nn.Dropout(config.attention_dropout)
18        self.resid_dropout = nn.Dropout(config.residual_dropout)
19
20        self.apply(self._init_weights)
21
22    def _init_weights(self, module):
23        if isinstance(module, nn.Linear):
24            torch.nn.init.normal_(module.weight, mean=0.0, std=0.02)
25            if module.bias is not None:
26                torch.nn.init.zeros_(module.bias)
27
28

```

```

29     def forward(self, q,k,v):
30         b,t,c = q.shape
31
32         q = self.q(q)
33         k = self.k(k)
34         v = self.v(v)
35
36         q = q.view(b,q.size(1),self.n_heads,self.head_size).permute(0,2,1,3) # batch x n_heads x seq_len x head_dim
37         k = k.view(b,k.size(1),self.n_heads,self.head_size).permute(0,2,1,3)
38         v = v.view(b,v.size(1),self.n_heads,self.head_size).permute(0,2,1,3)
39
40         qk_t = (q@k.transpose(-2,-1)) * self.scale
41         qk_t = F.softmax(qk_t,dim=-1)
42         weights = self.attn_dropout(qk_t)
43
44         attention = weights @ v # batch x n_heads x t x head_size
45         attention = attention.permute(0,2,1,3).contiguous().view(b,t,c) # batch x t x embed_dim
46
47         out = self.c_proj(attention)
48         out = self.resid_dropout(out)
49
50     return out

```

Feed Forward Block

```

1 class GPT2MLP(nn.Module):
2     def __init__(self,config):
3         super().__init__()
4         self.embed_dim = config.embed_dim
5         self.mlp_ratio = config.mlp_ratio
6         self.mlp_dropout = config.mlp_dropout
7
8         self.c_fc = nn.Linear(self.embed_dim, self.embed_dim*self.mlp_ratio)
9         self.c_proj = nn.Linear(self.embed_dim*self.mlp_ratio, self.embed_dim)
10        self.act = nn.GELU()
11        self.dropout = nn.Dropout(self.mlp_dropout)
12
13    def forward(self,x):
14        x = self.c_fc(x)
15        x = self.act(x)
16        x = self.c_proj(x)
17        x = self.dropout(x)
18        return x

```

Decoder Block

- with added cross-attention and pre-normalization

```

1 class GPT2Block(nn.Module):
2     def __init__(self,config):
3         super().__init__()
4         self.embed_dim = config.embed_dim
5         self.ln_1 = nn.LayerNorm(self.embed_dim)
6         self.attn = GPT2Attention(config)
7         self.ln_2 = nn.LayerNorm(self.embed_dim)
8         self.mlp = GPT2MLP(config)
9         self.ln_3 = nn.LayerNorm(self.embed_dim)
10        self.cross_attn = GPT2CrossAttention(config)
11
12    def forward(self,x,enc_out):
13        x = x+self.attn(self.ln_1(x))
14        x = x+self.cross_attn(self.ln_2(x),enc_out,enc_out)
15        x = x+self.mlp(self.ln_3(x))
16        return x

```

The main model

- creating ViT model via timm
- added embedding code and right forward pass for the ViT model as per [their code](#).
- each encoder output is passed to the decoder input and to the next encoder input as well.
- the final decoder outputs are passed through a head block to generate logits as per the vocab size.
- loss is calculated using cross-entropy if labels are present else, logits for the final token in the sequence are returned for generation.

- `pretrained_layers_trainable`: freezes/unfreezes ViT and GPT2 pretrained layers
- `unfreeze_gpt_layers`: unfreezes GPT2 layers only
- `from_pretrained`: loads GPT2 weights via huggingface gpt2
- `generate`: generates caption, sampling via `torch.multinomial` or deterministic via `argmax` with temperature control.

```

1 class VisionGPT2Model(nn.Module):
2     def __init__(self, config):
3         super().__init__()
4
5         self.config = config
6
7         vit = create_model('vit_base_patch16_224', pretrained=True, num_classes=0)
8         self.patch_embed = vit.patch_embed
9         num_patches = self.patch_embed.num_patches
10
11         self.cls_token = vit.cls_token
12         embed_len = num_patches + vit.num_prefix_tokens
13         self.pos_embed = vit.pos_embed
14         self.pos_drop = nn.Dropout(p=0.)
15
16         self.blocks = nn.ModuleList([vit.blocks[i] for i in range(config.depth)])
17
18         self.transformer = nn.ModuleDict(dict(
19             wte = nn.Embedding(config.vocab_size, config.embed_dim),
20             wpe = nn.Embedding(config.seq_len, config.embed_dim),
21             drop = nn.Dropout(config.emb_dropout),
22             h = nn.ModuleList([GPT2Block(config) for _ in range(config.depth)]),
23             ln_f = nn.LayerNorm(config.embed_dim)
24         ))
25         self.lm_head = nn.Linear(config.embed_dim, config.vocab_size, bias=False)
26         self.transformer.wte.weight = self.lm_head.weight
27
28     def _pos_embed(self, x):
29         pos_embed = self.pos_embed
30         x = torch.cat((self.cls_token.expand(x.shape[0], -1, -1), x), dim=1)
31         x = x + pos_embed
32         return self.pos_drop(x)
33
34     def pretrained_layers_trainable(self, trainable=False):
35         layers = [
36             self.cls_token, self.patch_embed, self.pos_embed, self.blocks,
37             self.transformer.wte, self.transformer.wpe,
38             self.transformer.ln_f, self.lm_head
39         ]
40         gpt_layers = [[
41             self.transformer.h[i].ln_1, self.transformer.h[i].ln_2,
42             self.transformer.h[i].attn, self.transformer.h[i].mlp
43         ] for i in range(self.config.depth)]
44         for l in gpt_layers:
45             layers.extend(l)
46
47         for layer in layers:
48             if not isinstance(layer, nn.Parameter):
49                 for p in layer.parameters():
50                     p.requires_grad = trainable
51             else:
52                 layer.requires_grad = trainable
53
54         total_frozen_params = sum([p.numel() for p in self.parameters() if not p.requires_grad])
55         print(f'total_frozen_params={total_frozen_params}')
56
57     def unfreeze_gpt_layers(self,):
58         gpt_layers = [[
59             self.transformer.h[i].ln_1, self.transformer.h[i].ln_2,
60             self.transformer.h[i].attn, self.transformer.h[i].mlp
61         ] for i in range(self.config.depth)]
62         flatten = []
63         for l in gpt_layers:
64             flatten.extend(l)
65
66         for layer in flatten:
67             if not isinstance(layer, nn.Parameter):
68                 for p in layer.parameters():
69                     p.requires_grad = True
70             else:
71                 layer.requires_grad = True
72
73     @classmethod
74     def from_pretrained(self, config):
75         model = VisionGPT2Model(config)

```

```

76     sd = model.state_dict()
77     keys = sd.keys()
78     ignore_matches = ['blocks.', 'cross_attn.', 'ln_3', 'cls_token', 'pos_embed', 'patch_embed.', '.attn.mask']
79     vit_keys = [key for key in keys if any(match in key for match in ignore_matches)]
80     gpt_keys = [key for key in keys if key not in vit_keys]
81
82     gpt2_small = GPT2LMHeadModel.from_pretrained('gpt2')
83     sd_hf = gpt2_small.state_dict()
84     hf_keys = sd_hf.keys()
85     hf_keys = [k for k in hf_keys if not k.endswith('.attn.masked_bias')]
86     hf_keys = [k for k in hf_keys if not k.endswith('.attn.bias')]
87     transposed = ['attn.c_attn.weight', 'attn.c_proj.weight', 'mlp.c_fc.weight', 'mlp.c_proj.weight']
88
89     for k in hf_keys:
90         if any(match in k for match in ignore_matches):
91             continue
92         if any(k.endswith(w) for w in transposed):
93             assert sd_hf[k].shape[::-1] == sd[k].shape
94             with torch.no_grad():
95                 sd[k].copy_(sd_hf[k].t())
96         else:
97             assert sd_hf[k].shape == sd[k].shape
98             with torch.no_grad():
99                 sd[k].copy_(sd_hf[k])
100
101     model.load_state_dict(sd)
102
103     return model
104
105 def forward(self,image,input_ids,labels=None):
106
107     image = self.patch_embed(image)
108     image = self._pos_embed(image)
109
110     token_embeddings = self.transformer.wte(input_ids) # batch x seq_len
111     pos_embs = torch.arange(0,input_ids.size(1)).to(input_ids.device)
112     positional_embeddings = self.transformer.wpe(pos_embs)
113     input_ids = self.transformer.drop(token_embeddings+positional_embeddings)
114
115     for i in range(self.config.depth):
116         image = self.blocks[i](image)
117         input_ids = self.transformer.h[i](input_ids, image)
118
119     input_ids = self.transformer.ln_f(input_ids)
120
121     if labels is not None:
122         lm_logits = self.lm_head(input_ids)
123         loss = F.cross_entropy(lm_logits.view(-1, lm_logits.shape[-1]), labels.view(-1))
124         return loss
125
126     lm_logits = self.lm_head(input_ids[:, [-1], :])
127     return lm_logits
128
129 def generate(self,image,sequence,max_tokens=50,temperature=1.0,deterministic=False):
130     for _ in range(max_tokens):
131         out = self(image,sequence)
132         out = out[:, :-1, :] / temperature
133         probs = F.softmax(out,dim=-1)
134         if deterministic:
135             next_token = torch.argmax(probs, dim=-1, keepdim=True)
136         else:
137             next_token = torch.multinomial(probs, num_samples=1)
138         sequence = torch.cat([sequence, next_token], dim=1)
139         if next_token.item() == tokenizer.eos_token_id:
140             break
141
142     return sequence.cpu().flatten()

```

▼ Training

- the pretrained layers are initially frozen as I need to train the cross attention layers first
- in the following epochs, GPT2 is unfreezed and trained, in the final few epochs, the ViT is unfreezed as well.
- optimizer: Adam
- scheduler: OneCycleLR
- mixed-precision fp16 training with autocast and grad-scaler in torch
- metrics: cross-entropy loss and perplexity = e^{loss} , both lower the better
- best model is saved based on validation perplexity and the same is loaded while generating captions.

Generation

- GPT2 generally requires context before generating anything, since for image captioning we can't really provide an initial context except the image itself, the initial context we provide to GPT is just [50256] i.e <|endoftext|> which is also the beginning of sentence token in GPT. For other models such as OPT it is </s>. This one token acts as the initial context.

```

1 class Trainer:
2     def __init__(self,model_config,train_config, dls):
3
4         self.train_config = train_config
5         self.model_config = model_config
6         self.device = self.train_config.device
7
8         self.model = VisionGPT2Model.from_pretrained(model_config).to(self.device)
9         self.model.pretrained_layers_trainable(trainable=False)
10
11     print(f'trainable parameters: {sum([p.numel() for p in self.model.parameters() if p.requires_grad])}')
12
13     self.tokenizer = GPT2TokenizerFast.from_pretrained('gpt2')
14     self.tokenizer.pad_token = self.tokenizer.eos_token
15
16     self.scaler = GradScaler()
17
18     self.train_dl, self.val_dl = dls
19
20     total_steps = len(self.train_dl)
21
22     self.optim = torch.optim.Adam(self.model.parameters(), lr=self.train_config.lr / 25.)
23     self.sched = torch.optim.lr_scheduler.OneCycleLR(
24         self.optim,
25         max_lr=self.train_config.lr,
26         epochs=self.train_config.epochs,
27         steps_per_epoch=total_steps
28     )
29
30 #     self.sched = get_linear_schedule_with_warmup(self.optim,num_warmup_steps=0,num_training_steps=total_steps)
31
32     self.metrics = pd.DataFrame()
33     self.metrics[['train_loss','train_perplexity','val_loss','val_perplexity']] = None
34
35     self.gen_tfms = A.Compose([
36         A.Resize(224,224),
37         A.Normalize(mean=[0.5,0.5,0.5],std=[0.5,0.5,0.5],always_apply=True),
38         ToTensorV2()
39     ])
40
41
42     def save_model(self,):
43         self.train_config.model_path.mkdir(exist_ok=True)
44         sd = self.model.state_dict()
45         torch.save(sd,self.train_config.model_path/'captioner.pt')
46
47
48     def load_best_model(self,):
49         sd = torch.load(self.train_config.model_path/'captioner.pt')
50         self.model.load_state_dict(sd)
51
52
53     def train_one_epoch(self,epoch):
54
55         prog = tqdm(self.train_dl,total=len(self.train_dl))
56
57         running_loss = 0.
58
59         for image, input_ids, labels in prog:
60
61             with autocast():
62                 image = image.to(self.device)
63                 input_ids = input_ids.to(self.device)
64                 labels = labels.to(self.device)
65
66                 loss = self.model(image,input_ids,labels)
67
68                 self.scaler.scale(loss).backward()
69                 self.scaler.step(self.optim)
70                 self.scaler.update()
71                 self.sched.step()
72                 self.optim.zero_grad(set_to_none=True)
73
74                 running_loss += loss.item()

```

```
75
76     prog.set_description(f'train loss: {loss.item():.3f}')
77
78     del image, input_ids, labels, loss
79
80     train_loss = running_loss / len(self.train_dl)
81     train_pxp = np.exp(train_loss)
82
83     self.metrics.loc[epoch,['train_loss','train_perplexity']] = (train_loss,train_pxp)
84
85
86     @torch.no_grad()
87     def valid_one_epoch(self,epoch):
88
89         prog = tqdm(self.val_dl,total=len(self.val_dl))
90
91         running_loss = 0.
92
93         for image, input_ids, labels in prog:
94
95             with autocast():
96                 image = image.to(self.device)
97                 input_ids = input_ids.to(self.device)
98                 labels = labels.to(self.device)
99
100            loss = self.model(image,input_ids,labels)
101            running_loss += loss.item()
102
103            prog.set_description(f'valid loss: {loss.item():.3f}')
104
105            del image, input_ids, labels, loss
106
107        val_loss = running_loss / len(self.val_dl)
108        val_pxp = np.exp(val_loss)
109
110        self.metrics.loc[epoch,['val_loss','val_perplexity']] = (val_loss,val_pxp)
111
112        return val_pxp
113
114
115    def clean(self):
116        gc.collect()
117        torch.cuda.empty_cache()
118
119
120    def fit(self,):
121
122        best_pxp = 1e9
123        best_epoch = -1
124        prog = tqdm(range(self.train_config.epochs))
125
126        for epoch in prog:
127
128            if epoch == self.train_config.freeze_epochs_gpt:
129                self.model.unfreeze_gpt_layers()
130                print('unfreezing GPT2 entirely...')
131
132            if epoch == self.train_config.freeze_epochs_all:
133                self.model.pretrained_layers_trainable(trainable=True)
134
135            self.model.train()
136            prog.set_description('training')
137            self.train_one_epoch(epoch)
138            self.clean()
139
140            self.model.eval()
141            prog.set_description('validating')
142            pxp = self.valid_one_epoch(epoch)
143            self.clean()
144
145            print(self.metrics.tail(1))
146
147            if pxp < best_pxp:
148                best_pxp = pxp
149                best_epoch = epoch
150                print('saving best model...')
151                self.save_model()
152
153        return {
154            'best_perplexity': best_pxp,
155            'best_epoch': best_epoch
156        }
```

```
157
158
159     @torch.no_grad()
160     def generate_caption(self,image,max_tokens=50,temperature=1.0,deterministic=False):
161         self.model.eval()
162
163         image = Image.open(image).convert('RGB')
164         image = np.array(image)
165         image = self.gen_tfms(image=image)['image']
166         image = image.unsqueeze(0).to(self.device)
167         sequence = torch.ones(1,1).to(device=self.device).long() * self.tokenizer.bos_token_id
168
169         caption = self.model.generate(
170             image,
171             sequence,
172             max_tokens=max_tokens,
173             temperature=temperature,
174             deterministic=deterministic
175         )
176         caption = self.tokenizer.decode(caption.numpy(),skip_special_tokens=True)
177
178
179     return caption

1 model_config = SimpleNamespace(
2     vocab_size = 50_257,
3     embed_dim = 768,
4     num_heads = 12,
5     seq_len = 1024,
6     depth = 12,
7     attention_dropout = 0.1,
8     residual_dropout = 0.1,
9     mlp_ratio = 4,
10    mlp_dropout = 0.1,
11    emb_dropout = 0.1,
12 )
13 train_config = SimpleNamespace(
14     epochs = 2,
15     freeze_epochs_gpt = 1,
16     freeze_epochs_all = 2,
17     lr = 1e-4,
18     device = 'cuda',
19     model_path = Path('captioner'),
20     batch_size = 32
21 )

1 train_dl = torch.utils.data.DataLoader(train_ds,batch_size=train_config.batch_size,shuffle=True,pin_memory=True,num_workers=2,persistent_workers=True)
2 val_dl = torch.utils.data.DataLoader(val_ds,batch_size=train_config.batch_size,shuffle=False,pin_memory=True,num_workers=2,persistent_workers=True)

1 trainer = Trainer(model_config,train_config,(train_dl,val_dl))

⤵ total_frozen_params=210236928
  trainable parameters: 28366848
  /usr/local/lib/python3.10/dist-packages/transformers/tokenization_utils_base.py:1601: FutureWarning: `clean_up_tokenization_spaces` is deprecated. Please use `clean_up_tokenization_spaces` instead.
  warnings.warn(
<ipython-input-28-4e526974c5c0>:16: FutureWarning: `torch.cuda.amp.GradScaler(args...)` is deprecated. Please use `torch.amp.GradScaler` instead.
  self.scaler = GradScaler()

◀ ➤

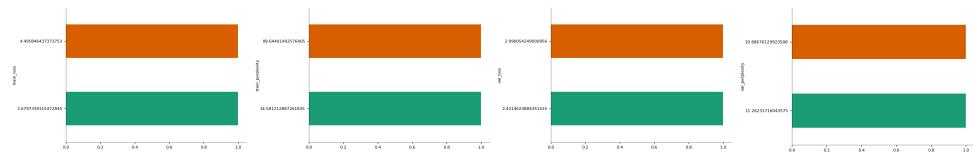
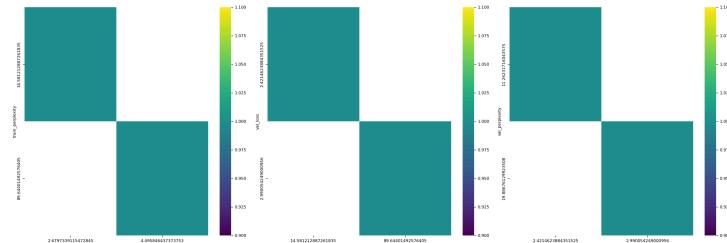
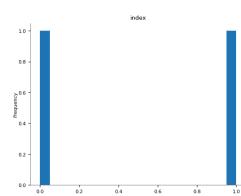
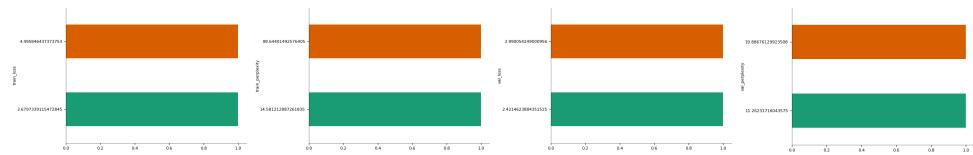
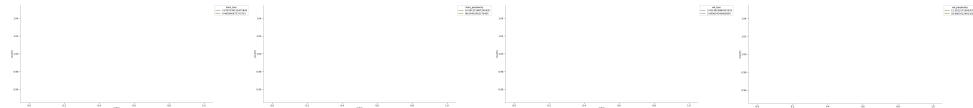
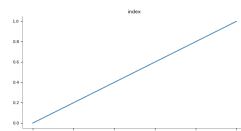
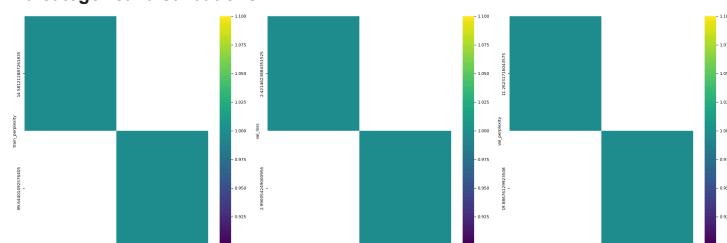
1 trainer.fit()
```

```
validate validating: 100%  
train loss: 3.357: 100% 2/2 [1:58:57<00:00, 3584.01s/it]  
<ipython-input-28-4e526974c5c0>:61: FutureWarning: `torch.cuda.amp.autocast(args...)` is deprecated. Please use `torch.amp.autocast with autocast():`  
valid loss: 2.593: 100% 4219/4219 [55:47<00:00, 1.39it/s]  
<ipython-input-28-4e526974c5c0>:95: FutureWarning: `torch.cuda.amp.autocast(args...)` is deprecated. Please use `torch.amp.autocast with autocast():`  
train_loss train_perplexity val_loss val_perplexity 469/469 [02:08<00:00, 5.01it/s]  
0 4.495846 89.644015 2.990054 19.886761  
saving best model...  
unfreezing GPT2 entirely...  
train loss: 2.847: 100% 4219/4219 [58:39<00:00, 2.05it/s]  
<ipython-input-28-4e526974c5c0>:61: FutureWarning: `torch.cuda.amp.autocast(args...)` is deprecated. Please use `torch.amp.autocast with autocast():`  
valid loss: 2.199: 100% 469/469 [02:04<00:00, 3.35it/s]  
<ipython-input-28-4e526974c5c0>:95: FutureWarning: `torch.cuda.amp.autocast(args...)` is deprecated. Please use `torch.amp.autocast with autocast():`  
train_loss train_perplexity val_loss val_perplexity 1 2.679734 14.581213 2.421462 11.262317  
saving best model...
```

▼ Results

```
1 trainer.metrics
```

	train_loss	train_perplexity	val_loss	val_perplexity
0	4.495846	89.644015	2.990054	19.886761
1	2.679734	14.581213	2.421462	11.262317

Categorical distributions**2-d categorical distributions****Distributions****Categorical distributions****Time series****Values****2-d categorical distributions****Faceted distributions**

<string>:5: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `le

<string>:5: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `le

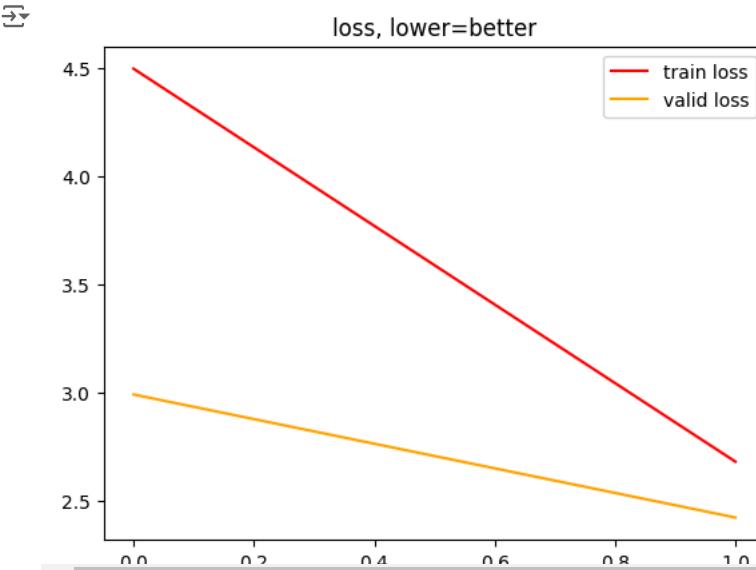
<string>:5: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `le

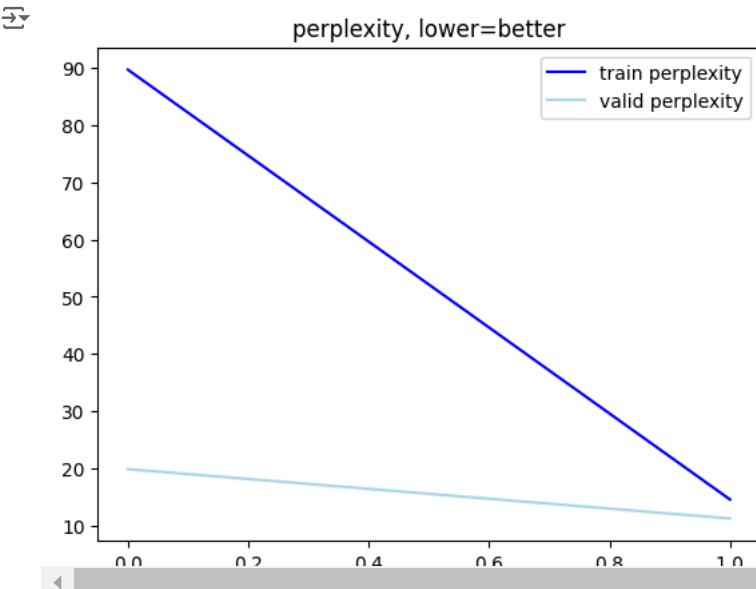
<string>:5: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `le

```
1 plt.plot(trainer.metrics['train_loss'],color='red',label='train loss')
2 plt.plot(trainer.metrics['val_loss'],color='orange',label='valid loss')
3 plt.title('loss, lower=better')
4 plt.legend()
5 plt.show()
```



```
1 plt.plot(trainer.metrics['train_perplexity'],color='blue',label='train perplexity')
2 plt.plot(trainer.metrics['val_perplexity'],color='lightblue',label='valid perplexity')
3 plt.title('perplexity, lower=better')
4 plt.legend()
5 plt.show()
```



▼ Predictions

```
1 trainer.load_best_model()
2 <ipython-input-28-4e526974c5c0>:49: FutureWarning: You are using `torch.load` with `weights_only=False` (the current default value),
3     sd = torch.load(self.train_config.model_path/'captioner.pt')

1 for i in range(50):
2     det = False
3     test = val_df.sample(n=1).values[0]
4     test_img, test_caption = test[0], test[1]
5     plt.imshow(Image.open(test_img).convert('RGB'))
6     t = np.random.uniform(0.5,1.5)
7     if i > 40:
```

```
8         det = True
9     gen_caption = trainer.generate_caption(test_img,temperature=t,deterministic=det)
10    plt.title(f"actual: {test_caption}\nmodel: {gen_caption}\ntemp: {t} deterministic generation: {det}")
11    plt.axis('off')
12    plt.show()
```



actual: A large very dimly lit kitchen with white cabinets.
 model: kitchen with a stove and oven in the kitchen.
 temp: 0.6425681266345307 deterministic generation: False



actual: A pile of banana sitting next to a stack of books.
 model: = some apples, bananas, bananas, and fruit.
 temp: 0.8758196029979488 deterministic generation: False



actual: A woman sitting at a table on a laptop.
 model: one of the benchescharging their burgers.
 temp: 1.363645290152989 deterministic generation: False



actual: The people are outside in the sun with umbrelalla over them
 model: and women walking in a field of people.
 temp: 0.6228119971393166 deterministic generation: False





actual: Rail yard with many tank and rail cars in large city.
 model: live trains waiting in the railroad yard a road and a black station.
 temp: 1.007670069937599 deterministic generation: False



actual: this plate has two pieces of bread on it
 model: wrap every sandwich on a blue plate featuring handwritten teddy bears.
 temp: 1.316646845769678 deterministic generation: False



actual: A vintage turquoise and beige truck is shown.
 model: top of a blue truck with the name on it
 temp: 0.6975312843490961 deterministic generation: False



actual: A woman in a bikini stands next to a surfboard.
 model: is a woman holding an unidentified surfboard in the ocean

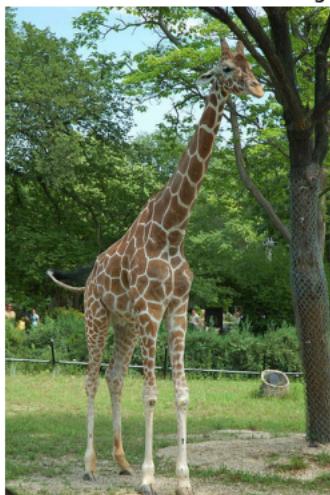
temp: 1.0293469299888411 deterministic generation: False



actual: A giraffe standing at attention in a zoo pin.

model: giraffes are standing at a zoo.

temp: 0.5238897311195714 deterministic generation: False



actual: A person stands by a wire fence as a cow looks around.
model: USE A BEHARISTIC ON A WHITE BOW WITH UPPIC ON THE SIDE

temp: 1.0044229990840443 deterministic generation: False



actual: A person in a yellow shirt swinging a tennis racket on a tennis court.
model: tennis player in action on court with racket.

temp: 0.564683112325976 deterministic generation: False

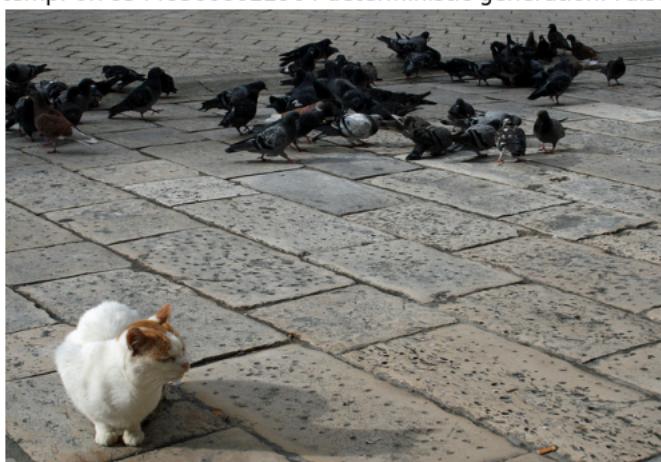




actual: A woman holds a plate of food while she has phone to her ear.
 model: closeup for picture showing two adequate tools a hand adapter someones higher changes cases
 temp: 1.4406396096151037 deterministic generation: False



actual: A white and orange cat in front of a flock of birds.
 model: irds standing in the shade near a large tree.
 temp: 0.7854483608022964 deterministic generation: False



actual: an elephant standing and looking into the camera
 model: dead tiring area with ape, stone tree impable meenary.
 temp: 1.4698097606799791 deterministic generation: False



actual: A table with a couch and joysticks on the table
 model: sitting in a couch beside a TV with a dog.
 temp: 0.8938652374742175 deterministic generation: False