```python
In [1]:  1  from models import Decoder, Img2Caption
         2
         3  import torch
         4  import torch.nn as nn
         5
         6  import torchvision
         7  from torchvision import transforms
         8
         9  import cv2
        10  from transformers import DebertaTokenizer, DebertaModel
```

```python
In [2]:  1  device = "cuda"
```

```python
In [3]:  1  tokenizer = DebertaTokenizer.from_pretrained("microsoft/deberta-base")
         2
         3  encoder_dim = 512
         4  encoder = torchvision.models.resnet101(pretrained=True)
         5  encoder.fc = nn.Linear(in_features=2048, out_features=encoder_dim)
         6
         7  decoder = Decoder(tokenizer.vocab_size, 768, encoder_dim, 256, 256, 0.2).to(device)
         8  gptModel = DebertaModel.from_pretrained('microsoft/deberta-base').to(device)
         9
        10  model = Img2Caption(encoder, decoder, gptModel).to(device)
```

Some weights of the model checkpoint at microsoft/deberta-base were not used when initializing DebertaModel: ['deberta.embeddings.position_embeddings.weight']
- This IS expected if you are initializing DebertaModel from the checkpoint of a model trained on another task or with another architecture (e.g. initializing a BertForSequenceClassification model from a BertForPreTraining model).
- This IS NOT expected if you are initializing DebertaModel from the checkpoint of a model that you expect to be exactly identical (initializing a BertForSequenceClassification model from a BertForSequenceClassification model).

```python
In [4]:  1  model.load_state_dict(torch.load("best_val_model.pt"))
```

Out[4]: <All keys matched successfully>

```python
In [5]:  1  def transform(image, dsize=(256, 256)):
         2      resized_image = cv2.resize(image, dsize)
         3      torch_transform = transforms.Compose(
         4          [transforms.ToTensor(), transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225])]
         5      )
         6      return torch_transform(resized_image)
         7
         8  def get_image(file_name):
         9      image = cv2.cvtColor(cv2.imread(file_name), cv2.COLOR_BGR2RGB)
        10      image = transform(image)
        11      return image
```

```python
In [6]:  1  def create_caption_from_image(image_path, model, max_caption_len):
         2  #     setting fake caption to init generation
         3      max_len = max_caption_len
         4      targets = torch.zeros(1, max_len).long().to(device)
         5      # 'cause tokenizer.encode("") results in [1, 2]
         6      targets[0, 1] = 1
         7
         8  #     loading and transforming image
         9      image = cv2.cvtColor(cv2.imread(image_path), cv2.COLOR_BGR2RGB)
        10      image = transform(image).unsqueeze(0).to(device)
        11
        12  #     making caption
        13      caption = model(image, targets, teacher_forcing_ratio=1)
        14      value = tokenizer.decode(torch.argmax(caption.permute(1, 0, 2), axis=2)[0])
        15
        16      return value
```

```
In [10]:    1  # captioning 1
            2  create_caption_from_image("kot.jpg", model, 32)
```

Out[10]: 'A woman with a a cat a a.. a'

```
In [8]:     1  # captiononig 2
            2  create_caption_from_image("kot.jpg", model, 32)
```

Out[8]: 'A small cat in a a a a a..'

```
In [ ]:     1
```