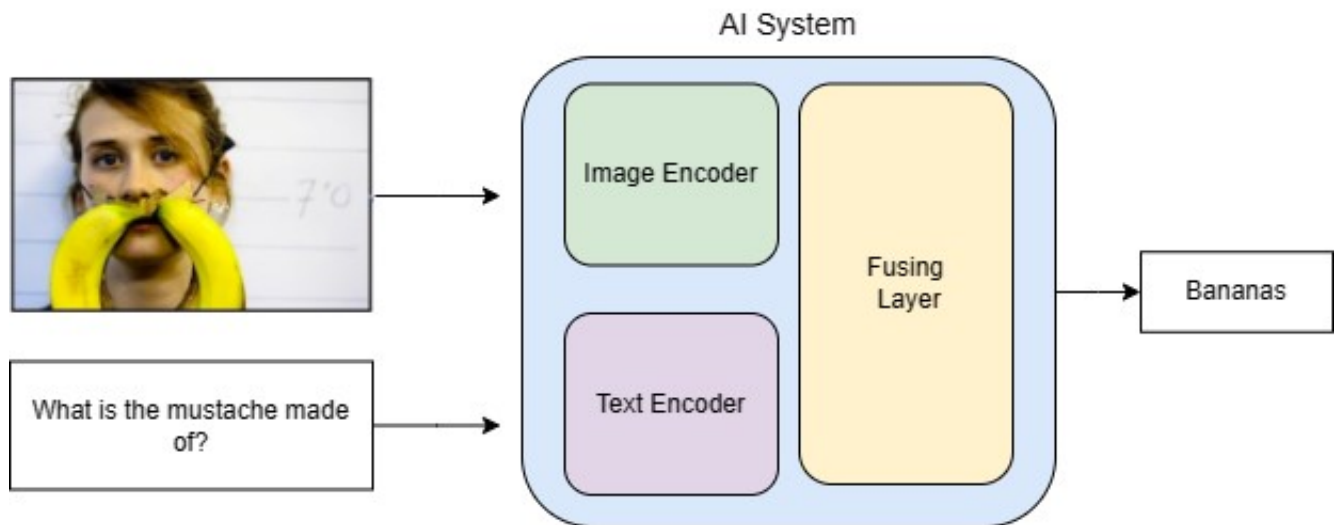


Visual Question Answering - Project

Ngày 10 tháng 12 năm 2023

Phần I: Giới thiệu

Visual Question Answering (VQA) là một bài toán phổ biến trong Machine Learning, ứng dụng các kĩ thuật liên quan từ hai lĩnh vực Computer Vision và Natural Language Processing. Khái niệm cốt lõi của bài toán này là phân tích một hình ảnh và trả lời câu hỏi về hình ảnh đó. Bước đầu là phân tích thông tin đầu vào, bao gồm sử dụng các kĩ thuật xử lý hình ảnh và xử lý câu hỏi đặt ra bằng ngôn ngữ tự nhiên. Sau đó, hệ thống VQA sẽ hợp thông tin thu được từ phân tích hình ảnh và ngữ cảnh của câu hỏi để tạo ra một câu trả lời phù hợp. Vì vậy, một chương trình có độ chính xác cao cần xây dựng tốt cả hai thành phần này, đặt ra thách thức rất lớn trong việc giải quyết tốt bài toán hỏi đáp với ảnh.



Trong project này, chúng ta sẽ xây dựng một chương trình VQA sử dụng mô hình CNN cho hình ảnh và LSTM cho xử lý ngôn ngữ. Input và output của chương trình như sau:

- **Input:** Một cặp hình ảnh và câu hỏi bằng ngôn ngữ tự nhiên.
- **Output:** Câu trả lời cho câu hỏi về hình ảnh.

Phần II: Cài đặt chương trình

A. Phần lập trình

- **CNN + LSTM:** Trong phần này, chúng ta sẽ tiếp cận bài toán bằng cách phối hợp hai mô hình Deep Learning cơ bản dùng trong xử lý ảnh và văn bản là mạng CNN và LSTM.

1. **Download dataset:** Cho bộ dữ liệu về Visual Question Answering có nội dung quan đến câu hỏi đáp dạng Yes/No, các bạn tải bộ dữ liệu này tại [đây](#). Một số mẫu từ bộ dữ liệu khi được trực quan hóa sẽ có dạng như hình bên dưới:



Hình 1: Một vài mẫu dữ liệu trong bộ dữ liệu VQA dạng câu hỏi Yes/No

2. **Import các thư viện cần thiết:** Chúng ta sẽ sử dụng thư viện PyTorch để xây dựng và huấn luyện mô hình deep learning. Thêm vào đó, vì làm việc liên quan đến dữ liệu ảnh và text, chúng ta sẽ sử dụng thư viện PIL cho ảnh và spacy, nltk cho text:

```
1 import torch
2 import torch.nn as nn
3 import os
4 import pandas as pd
5 import matplotlib.pyplot as plt
6
7 from PIL import Image
8 from torch.utils.data import Dataset, DataLoader
9 from sklearn.model_selection import train_test_split
10 import torchtext
11 from torchtext.data.utils import get_tokenizer
```

```
12 from torchtext.vocab import build_vocab_from_iterator
13 import spacy
14
```

3. **Chia bộ dữ liệu train, val, test:** Vì bộ dữ liệu này đã được chia sẵn thành train, val, test, ta chỉ cần đọc dữ liệu lên từ file .txt cho trước:

```
1 # Load train data
2 train_data = []
3 train_set_path = './vaq2.0.TrainImages.txt'
4
5 with open(train_set_path, "r") as f:
6     lines = f.readlines()
7     for line in lines:
8         temp = line.split('\t')
9         qa = temp[1].split('?')
10
11         if len(qa) == 3:
12             answer = qa[2].strip()
13         else:
14             answer = qa[1].strip()
15
16         data_sample = {
17             'image_path': temp[0][: -2],
18             'question': qa[0] + '?',
19             'answer': answer
20         }
21         train_data.append(data_sample)
22
23 # Load val data
24 val_data = []
25 val_set_path = './vaq2.0.DevImages.txt'
26
27 with open(val_set_path, "r") as f:
28     lines = f.readlines()
29     for line in lines:
30         temp = line.split('\t')
31         qa = temp[1].split('?')
32
33         if len(qa) == 3:
34             answer = qa[2].strip()
35         else:
36             answer = qa[1].strip()
37
38         data_sample = {
39             'image_path': temp[0][: -2],
40             'question': qa[0] + '?',
41             'answer': answer
42         }
43         val_data.append(data_sample)
44
45 # Load test data
46 test_data = []
47 test_set_path = './vaq2.0.TestImages.txt'
48
49 with open(test_set_path, "r") as f:
50     lines = f.readlines()
51     for line in lines:
52         temp = line.split('\t')
53         qa = temp[1].split('?')
54
```

```

55     if len(qa) == 3:
56         answer = qa[2].strip()
57     else:
58         answer = qa[1].strip()
59
60     data_sample = {
61         'image_path': temp[0][:-2],
62         'question': qa[0] + '?',
63         'answer': answer
64     }
65     test_data.append(data_sample)
66

```

4. **Xây dựng bộ từ vựng:** Chúng ta cần tiền xử lý text bằng cách biến đổi câu hỏi đầu vào thành các token bằng thư viện spacy và xây dựng bộ từ vựng cho model.

```

1  eng = spacy.load("en_core_web_sm") # Load the English model to tokenize
   English text
2
3  def get_tokens(data_iter):
4      for sample in data_iter:
5          question = sample['question']
6          yield [token.text for token in eng.tokenizer(question)]
7
8
9  vocab = build_vocab_from_iterator(
10     get_tokens(train_data),
11     min_freq=2,
12     specials= ['<pad>', '<sos>', '<eos>', '<unk>'],
13     special_first=True
14 )
15 vocab.set_default_index(vocab['<unk>'])
16

```

5. **Xây dựng dictionary mapping classes:** Để thuận tiện trong việc chuyển đổi tên class (trong trường hợp này gồm 2 class là "yes" và "no"), ta tạo dictionary dùng để chuyển tên class thành mã số tương ứng và ngược lại. Cách làm như sau:

```

1  classes = set([sample['answer'] for sample in train_data])
2  classes_to_idx = {
3      cls_name: idx for idx, cls_name in enumerate(classes)
4  }
5  idx_to_classes = {
6      idx: cls_name for idx, cls_name in enumerate(classes)
7  }
8

```

6. **Xây dựng hàm tokenize:** Sau khi đã có bộ từ vựng cho model, ta xây dựng một hàm tokenize để biến câu hỏi đầu vào thành danh sách các token tương ứng trong bộ từ vựng:

```

1  def tokenize(question, max_sequence_length):
2      tokens = [token.text for token in eng.tokenizer(question)]
3      sequence = [vocab[token] for token in tokens]
4      if len(sequence) < max_sequence_length:
5          sequence += [vocab['<pad>']] * (max_sequence_length - len(sequence))
6      else:
7          sequence = sequence[:max_sequence_length]
8
9      return sequence
10

```

7. **Xây dựng class pytorch dataset:** Chúng ta xây dựng class datasets cho bộ dữ liệu VQA như sau:

```

1 class VQADataset(Dataset):
2     def __init__(
3         self,
4         data,
5         classes_to_idx,
6         max_seq_len=30,
7         transform=None,
8         root_dir='/content/val2014-resised/'
9     ):
10         self.transform = transform
11         self.data = data
12         self.max_seq_len = max_seq_len
13         self.root_dir = root_dir
14         self.classes_to_idx = classes_to_idx
15
16     def __len__(self):
17         return len(self.data)
18
19     def __getitem__(self, index):
20         img_path = os.path.join(self.root_dir, self.data[index]['image_path'
21 ])
22         img = Image.open(img_path).convert('RGB')
23         if self.transform:
24             img = self.transform(img)
25
26         question = self.data[index]['question']
27         question = tokenize(question, self.max_seq_len)
28         question = torch.tensor(question, dtype=torch.long)
29
30         label = self.data[index]['answer']
31         label = classes_to_idx[label]
32         label = torch.tensor(label, dtype=torch.long)
33
34         return img, question, label

```

8. **Xây dựng hàm tiền xử lý ảnh (transforms):** Để dữ liệu ảnh đầu vào được đồng bộ về kích thước cũng như đơn giản hóa, chúng ta sẽ xây dựng hàm tiền xử lý ảnh như sau:

```

1 transform = transforms.Compose([
2     transforms.Resize((224, 224)),
3     transforms.ToTensor(),
4     transforms.Normalize((0.485, 0.456, 0.406), (0.229, 0.224, 0.225)),
5 ])
6

```

Các kỹ thuật được áp dụng: resize ảnh, đổi về tensor và chuẩn hóa giá trị pixel.

9. **Khai báo datasets object cho ba bộ train, val, test:**

```

1 train_dataset = VQADataset(
2     train_data,
3     classes_to_idx=classes_to_idx,
4     transform=transform
5 )
6 val_dataset = VQADataset(
7     val_data,
8     classes_to_idx=classes_to_idx,
9     transform=transform
10 )

```

```

11 test_dataset = VQADataset(
12     test_data,
13     classes_to_idx=classes_to_idx,
14     transform=transform
15 )

```

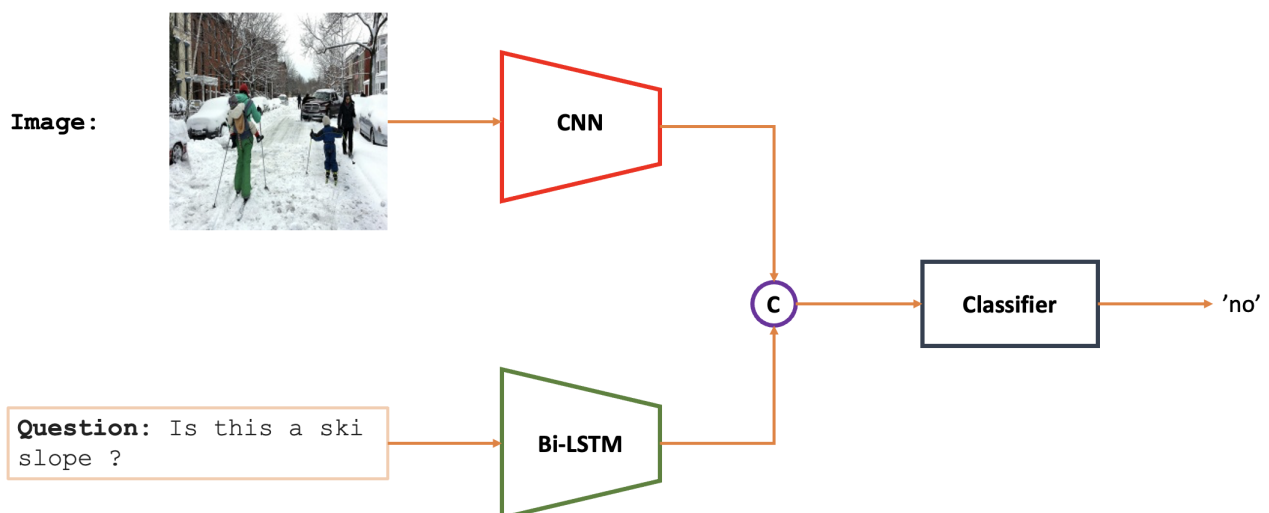
10. **Khai báo DataLoader:** Với ba object datasets trên, ta khai báo giá trị batch size và tạo dataloader như sau:

```

1 train_batch_size = 128
2 test_batch_size = 32
3
4 train_loader = DataLoader(
5     train_dataset,
6     batch_size=train_batch_size,
7     shuffle=True
8 )
9 val_loader = DataLoader(
10    val_dataset,
11    batch_size=test_batch_size,
12    shuffle=False
13 )
14 test_loader = DataLoader(
15    test_dataset,
16    batch_size=test_batch_size,
17    shuffle=False
18 )

```

11. **Xây dựng model:** Trong phần này, ta sẽ dùng kiến trúc ResNet cho phần xử lý ảnh và kiến trúc LSTM cho phần xử lý text. Hướng tiếp cận này có thể được mô phỏng qua ảnh sau:



Hình 2: Ảnh mô tả pipeline cho hướng tiếp cận sử dụng mô hình CNN để trích xuất đặc trưng ảnh và mô hình Bi-LSTM để trích xuất đặc trưng văn bản

Cụ thể, đối với ResNet, ta dùng thư viện [timm](#) để load kiến trúc ResNet50. Code của class model có nội dung như sau:

```

1 class VQAModel(nn.Module):
2     def __init__(
3         self,

```

```

4         n_classes,
5         img_model_name='resnet50',
6         embeddding_dim=300,
7         n_layers=2,
8         hidden_size=128,
9         dropout_prob=0.2
10    ):
11        super(VQAModel, self).__init__()
12        self.image_encoder = timm.create_model(
13            img_model_name,
14            pretrained=True,
15            num_classes=hidden_size
16        )
17
18        self.embedding = nn.Embedding(len(vocab), embeddding_dim)
19        self.lstm = nn.LSTM(
20            input_size=embeddding_dim,
21            hidden_size=hidden_size,
22            num_layers=n_layers,
23            batch_first=True,
24            bidirectional=True
25        )
26        self.layernorm = nn.LayerNorm(hidden_size * 2)
27        self.fc1 = nn.Linear(hidden_size * 3, 256)
28        self.relu = nn.ReLU()
29        self.dropout = nn.Dropout(dropout_prob)
30        self.fc2 = nn.Linear(256, n_classes)
31
32    def forward(self, img, text):
33        img_features = self.image_encoder(img)
34
35        text_emb = self.embedding(text)
36        lstm_out, _ = self.lstm(text_emb)
37
38        lstm_out = lstm_out[:, -1, :]
39        lstm_out = self.layernorm(lstm_out)
40
41        combined = torch.cat((img_features, lstm_out), dim=1)
42        x = self.fc1(combined)
43        x = self.relu(x)
44        x = self.dropout(x)
45        x = self.fc2(x)
46
47        return x
48
```

Sau khi định nghĩa class, ta tiến hành khai báo model như sau:

```

1 n_classes = len(classes)
2 img_model_name = 'resnet50'
3 hidden_size = 128
4 n_layers = 1
5 embeddding_dim = 128
6 dropout_prob = 0.2
7 device = 'cuda' if torch.cuda.is_available() else 'cpu'
8
9 model = VQAModel(
10     n_classes=n_classes,
11     img_model_name=img_model_name,
12     embeddding_dim=embeddding_dim,
13     n_layers=n_layers,

```

```

14     hidden_size=hidden_size,
15     dropout_prob=dropout_prob
16 ).to(device)
17

```

12. **Xây dựng hàm train và evaluate:** Để huấn luyện mô hình, ta cần xây dựng hàm huấn luyện và đánh giá mô hình như sau:

```

1  def evaluate(model, dataloader, criterion, device):
2      model.eval()
3      correct = 0
4      total = 0
5      losses = []
6      with torch.no_grad():
7          for image, question, labels in dataloader:
8              image, question, labels = image.to(device), question.to(device),
              labels.to(device)
9              outputs = model(image, question)
10             loss = criterion(outputs, labels)
11             losses.append(loss.item())
12             _, predicted = torch.max(outputs.data, 1)
13             total += labels.size(0)
14             correct += (predicted == labels).sum().item()
15
16     loss = sum(losses) / len(losses)
17     acc = correct / total
18
19     return loss, acc
20
21 def fit(
22     model,
23     train_loader,
24     val_loader,
25     criterion,
26     optimizer,
27     scheduler,
28     device,
29     epochs
30 ):
31     train_losses = []
32     val_losses = []
33
34     for epoch in range(epochs):
35         batch_train_losses = []
36
37         model.train()
38         for idx, (images, questions, labels) in enumerate(train_loader):
39             images = images.to(device)
40             questions = questions.to(device)
41             labels = labels.to(device)
42
43             optimizer.zero_grad()
44             outputs = model(images, questions)
45             loss = criterion(outputs, labels)
46             loss.backward()
47             optimizer.step()
48
49             batch_train_losses.append(loss.item())
50
51     train_loss = sum(batch_train_losses) / len(batch_train_losses)
52     train_losses.append(train_loss)

```



```

53         val_loss, val_acc = evaluate(
54             model, val_loader,
55             criterion, device
56         )
57         val_losses.append(val_loss)
58
59         print(f'EPOCH {epoch + 1}: \tTrain loss: {train_loss:.4f} \tVal loss: {
60             val_loss:.4f} \tVal Acc: {val_acc}')
61
62         scheduler.step()
63
64     return train_losses, val_losses
65

```

13. Khai báo hàm loss, optimizer và scheduler:

```

1  lr = 1e-2
2  epochs = 50
3  weight_decay = 1e-5
4  scheduler_step_size = epochs * 0.6
5  criterion = nn.CrossEntropyLoss()
6
7  optimizer = torch.optim.Adam(
8      model.parameters(),
9      lr=lr,
10     weight_decay=weight_decay
11 )
12 scheduler = torch.optim.lr_scheduler.StepLR(
13     optimizer,
14     step_size=scheduler_step_size,
15     gamma=0.1
16 )
17

```

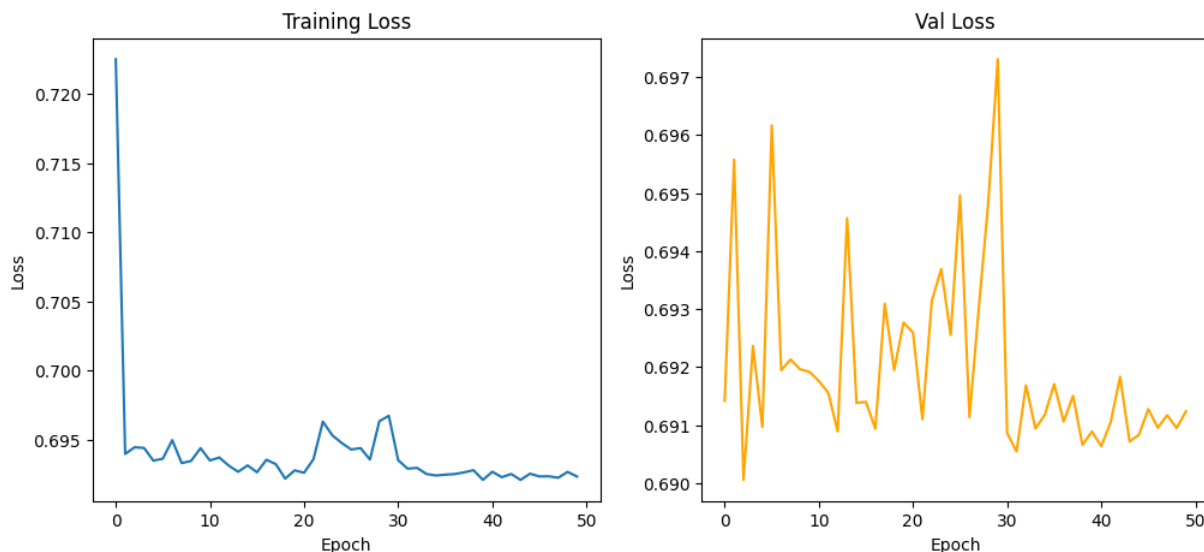
14. **Training:** Tổng hợp tất cả các thành phần trên, ta gọi hàm `fit()` để bắt đầu quá trình huấn luyện mô hình VQA:

```

1  train_losses, val_losses = fit(
2      model,
3      train_loader,
4      val_loader,
5      criterion,
6      optimizer,
7      scheduler,
8      device,
9      epochs
10 )
11

```

Khi quá trình huấn luyện kết thúc, ta có thể trực quan kết quả loss của mô hình qua từng epoch như sau:



Hình 3: Kết quả huấn luyện của mô hình với ResNet+Bi-LSTM

15. **Evaluation:** Với mô hình đã huấn luyện, ta đưa vào đánh giá trên hai tập val và test như sau:

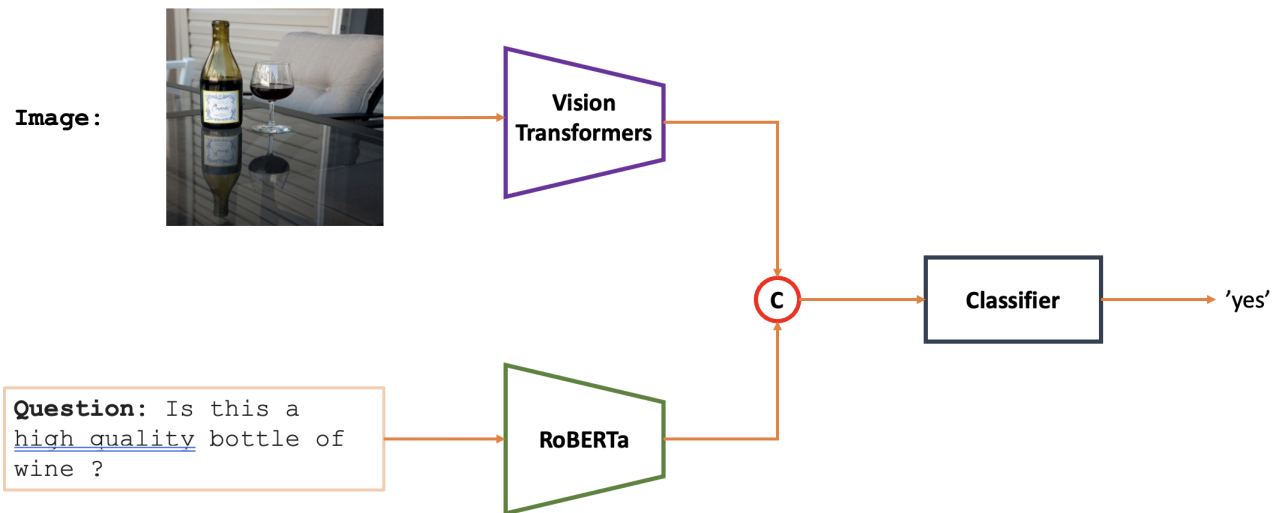
```

1 val_loss, val_acc = evaluate(
2     model,
3     val_loader,
4     criterion,
5     device
6 )
7 test_loss, test_acc = evaluate(
8     model,
9     test_loader,
10    criterion,
11    device
12 )
13
14 print('Evaluation on val/test dataset')
15 print('Val accuracy: ', val_acc)
16 print('Test accuracy: ', test_acc)
17

```

Có thể thấy, với hướng tiếp cận cơ bản trên, mô hình VQA ta xây dựng không đạt được kết quả khả quan. Vì vậy ở phần sau, ta sẽ sử dụng các mô hình tốt hơn để cải thiện kết quả bài toán này.

- **VisionTransformers + RoBERTa:** Trong phần này, chúng ta sẽ tiếp cận bài toán bằng cách phối hợp hai mô hình Deep Learning thuộc họ mô hình transformers, một kiến trúc mạng cực kì mạnh mẽ đang được sử dụng rộng rãi tại thời điểm này. Cụ thể, với ảnh đầu vào, ta dùng VisionTransformers, với câu hỏi, ta dùng RoBERTa. Với kết quả từ hai mô hình, ta sẽ kết hợp chúng để dự đoán câu trả lời. Hướng tiếp cận này có thể được mô phỏng như hình sau:



Hình 4: Ảnh mô tả pipeline cho hướng tiếp cận sử dụng các mô hình thuộc họ transformers

Chúng ta sẽ triển khai hướng tiếp cận này thông qua các bước thực hiện dưới đây:

1. Import các thư viện cần thiết:

```

1 import torch
2 import torch.nn as nn
3 import os
4 import numpy as np
5 import pandas as pd
6 import timm
7 import matplotlib.pyplot as plt
8
9 from PIL import Image
10 from torch.utils.data import Dataset, DataLoader
11 from transformers import ViTModel, ViTImageProcessor
12 from transformers import AutoTokenizer, RobertaModel
13

```

2. Chia bộ dữ liệu train, val, test: Tương tự phần trên, ta tạo ba bộ dữ liệu train, val, test:

```

1 # Load train data
2 train_data = []
3 train_set_path = './vaq2.0.TrainImages.txt'
4
5 with open(train_set_path, "r") as f:
6     lines = f.readlines()
7     for line in lines:
8         temp = line.split('\t')
9         qa = temp[1].split('?')
10
11         if len(qa) == 3:
12             answer = qa[2].strip()
13         else:
14             answer = qa[1].strip()
15
16         data_sample = {
17             'image_path': temp[0][:-2],

```

```

18         'question': qa[0] + '?',
19         'answer': answer
20     }
21     train_data.append(data_sample)
22
23 # Load val data
24 val_data = []
25 val_set_path = './vaq2.0.DevImages.txt'
26
27 with open(val_set_path, "r") as f:
28     lines = f.readlines()
29     for line in lines:
30         temp = line.split('\t')
31         qa = temp[1].split('?')
32
33         if len(qa) == 3:
34             answer = qa[2].strip()
35         else:
36             answer = qa[1].strip()
37
38         data_sample = {
39             'image_path': temp[0][:-2],
40             'question': qa[0] + '?',
41             'answer': answer
42         }
43         val_data.append(data_sample)
44
45 # Load test data
46 test_data = []
47 test_set_path = './vaq2.0.TestImages.txt'
48
49 with open(test_set_path, "r") as f:
50     lines = f.readlines()
51     for line in lines:
52         temp = line.split('\t')
53         qa = temp[1].split('?')
54
55         if len(qa) == 3:
56             answer = qa[2].strip()
57         else:
58             answer = qa[1].strip()
59
60         data_sample = {
61             'image_path': temp[0][:-2],
62             'question': qa[0] + '?',
63             'answer': answer
64         }
65         test_data.append(data_sample)
66

```

3. **Xây dựng dictionary mapping classes:** Tương tự như phần trên, ta khởi tạo dictionary dùng để chuyển đổi tên class thành mã số tương ứng và ngược lại:

```

1 classes = set([sample['answer'] for sample in train_data])
2 classes_to_idx = {
3     cls_name: idx for idx, cls_name in enumerate(classes)
4 }
5 idx_to_classes = {
6     idx: cls_name for idx, cls_name in enumerate(classes)
7 }
8

```

4. **Xây dựng class pytorch dataset:** Với việc sử dụng hai mô hình liên quan đến transformers, chúng ta sẽ có một chút sự thay đổi trong việc triển khai class dataset như sau:

```

1 class VQADataset(Dataset):
2     def __init__(
3         self,
4         data,
5         classes_to_idx,
6         img_feature_extractor,
7         text_tokenizer,
8         device,
9         root_dir='/content/val2014-resised/'
10    ):
11        self.data = data
12        self.root_dir = root_dir
13        self.classes_to_idx = classes_to_idx
14        self.img_feature_extractor = img_feature_extractor
15        self.text_tokenizer = text_tokenizer
16        self.device = device
17
18    def __len__(self):
19        return len(self.data)
20
21    def __getitem__(self, index):
22        img_path = os.path.join(self.root_dir, self.data[index]['image_path'
23    ])
24
25        img = Image.open(img_path).convert('RGB')
26
27        if self.img_feature_extractor:
28            img = self.img_feature_extractor(images=img, return_tensors="pt")
29            img = {k: v.to(self.device).squeeze(0) for k, v in img.items()}
30
31        question = self.data[index]['question']
32        if self.text_tokenizer:
33            question = self.text_tokenizer(
34                question,
35                padding="max_length",
36                max_length=20,
37                truncation=True,
38                return_tensors="pt"
39            )
40            question = {k: v.to(self.device).squeeze(0) for k, v in question.
41    items()}
42
43        label = self.data[index]['answer']
44        label = torch.tensor(
45            classes_to_idx[label],
46            dtype=torch.long
47        ).to(device)
48
49        sample = {
50            'image': img,
51            'question': question,
52            'label': label
53        }
54
55        return sample

```

Ở đây, ta bổ sung vào hai thành phần mới gồm `img_feature_extractor` và `text_tokenizer`, đây là

hai hàm dùng để tiền xử lý dữ liệu đầu vào dành riêng cho VisionTransformers và RoBERTa. Vì kết quả của hai hàm này là tensor, vì vậy chúng ta cũng cần truyền vào tham số `device` để chuyển đổi tensor về format tính toán của máy.

5. Khai báo datasets object cho ba bộ train, val, test:

```

1 img_feature_extractor = ViTImageProcessor.from_pretrained("google/vit-base-
  patch16-224")
2 text_tokenizer = AutoTokenizer.from_pretrained("roberta-base")
3 device = 'cuda' if torch.cuda.is_available() else 'cpu'
4
5 train_dataset = VQADataset(
6     train_data,
7     classes_to_idx=classes_to_idx,
8     img_feature_extractor=img_feature_extractor,
9     text_tokenizer=text_tokenizer,
10    label_encoder=label_encoder,
11    device=device
12 )
13 val_dataset = VQADataset(
14     val_data,
15     classes_to_idx=classes_to_idx,
16     img_feature_extractor=img_feature_extractor,
17     text_tokenizer=text_tokenizer,
18     label_encoder=label_encoder,
19     device=device
20 )
21 test_dataset = VQADataset(
22     test_data,
23     classes_to_idx=classes_to_idx,
24     img_feature_extractor=img_feature_extractor,
25     text_tokenizer=text_tokenizer,
26     label_encoder=label_encoder,
27     device=device
28 )
29

```

6. **Xây dựng model:** Ta sẽ chia mô hình thành 3 phần gồm TextEncoder, VisualEncoder và Classifier. Sau đó, ta kết hợp ba thành phần này lại để trở thành một mô hình hoàn chỉnh. Code triển khai như sau:

– TextEncoder:

```

1 class TextEncoder(nn.Module):
2     def __init__(self):
3         super(TextEncoder, self).__init__()
4         self.model = RobertaModel.from_pretrained("roberta-base")
5
6     def forward(self, inputs):
7         outputs = self.model(**inputs)
8
9         return outputs.pooler_output
10

```

– VisualEncoder:

```

1 class VisualEncoder(nn.Module):
2     def __init__(self):
3         super(VisualEncoder, self).__init__()
4         self.model = ViTModel.from_pretrained("google/vit-base-patch16-224")
5

```

```

6     def forward(self, inputs):
7         outputs = self.model(**inputs)
8
9         return outputs.pooler_output
10

```

— **Classifier:**

```

1 class Classifier(nn.Module):
2     def __init__(
3         self,
4         input_size=768*2,
5         hidden_size=512,
6         n_layers=1,
7         dropout_prob=0.2,
8         n_classes=2
9     ):
10        super(Classifier, self).__init__()
11        self.lstm = nn.LSTM(
12            input_size,
13            hidden_size,
14            num_layers=n_layers,
15            batch_first=True,
16            bidirectional=True
17        )
18        self.dropout = nn.Dropout(dropout_prob)
19        self.fc1 = nn.Linear(hidden_size*2, n_classes)
20
21    def forward(self, x):
22        x, _ = self.lstm(x)
23        x = self.dropout(x)
24        x = self.fc1(x)
25
26        return x
27

```

— **VQAModel:** Tổng hợp lại 3 thành phần trên, ta được mô hình VQA hoàn chỉnh như sau:

```

1 class VQAModel(nn.Module):
2     def __init__(
3         self,
4         visual_encoder,
5         text_encoder,
6         classifier
7     ):
8        super(VQAModel, self).__init__()
9        self.visual_encoder = visual_encoder
10       self.text_encoder = text_encoder
11       self.classifier = classifier
12
13    def forward(self, image, answer):
14        text_out = self.text_encoder(answer)
15        image_out = self.visual_encoder(image)
16        x = torch.cat((text_out, image_out), dim=1)
17        x = self.classifier(x)
18
19        return x
20
21    def freeze(self, visual=True, textual=True, clas=False):
22        if visual:
23            for n,p in self.visual_encoder.named_parameters():

```

```

24         p.requires_grad = False
25     if textual:
26         for n,p in self.text_encoder.named_parameters():
27             p.requires_grad = False
28     if clas:
29         for n,p in self.classifier.named_parameters():
30             p.requires_grad = False
31

```

Các bạn có thể thấy trong phần code này, chúng ta có triển khai thêm hàm `freeze()`. Hàm này có chức năng dùng để đóng băng các tham số của những pretrained model, cụ thể ở đây là một trong ba thành phần chính của mô hình. Trong bài này, chúng ta sẽ chỉ cập nhật trọng số cho phần Classifier, các phần khác chúng ta sẽ đóng băng lại.

Sau khi định nghĩa xong, ta khai báo mô hình VQA:

```

1 n_classes = len(classes)
2 hidden_size = 1024
3 n_layers = 1
4 dropout_prob = 0.2
5
6 text_encoder = TextEncoder().to(device)
7 visual_encoder = VisualEncoder().to(device)
8 classifier = Classifier(
9     hidden_size=hidden_size,
10    n_layers=n_layers,
11    dropout_prob=dropout_prob,
12    n_classes=n_classes
13).to(device)
14
15 model = VQAModel(
16     visual_encoder=visual_encoder,
17     text_encoder=text_encoder,
18     classifier=classifier
19).to(device)
20 model.freeze()
21

```

7. **Xây dựng hàm train và evaluate:** Ta xây dựng hàm train và evaluate. Ở phần này, nội dung code của hai hàm sẽ có đôi chút sự khác biệt trong việc lấy dữ liệu từ dataloader vì chúng ta đã thay đổi cấu trúc data lúc định nghĩa pytorch dataset:

```

1 def evaluate(model, dataloader, criterion):
2     model.eval()
3     correct = 0
4     total = 0
5     losses = []
6     with torch.no_grad():
7         for idx, inputs in enumerate(dataloader):
8             images = inputs['image']
9             questions = inputs['question']
10            labels = inputs['label']
11            outputs = model(images, questions)
12            loss = criterion(outputs, labels)
13            losses.append(loss.item())
14            _, predicted = torch.max(outputs.data, 1)
15            total += labels.size(0)
16            correct += (predicted == labels).sum().item()
17
18     loss = sum(losses) / len(losses)
19     acc = correct / total

```



```

20
21     return loss, acc
22
23 def fit(
24     model,
25     train_loader,
26     val_loader,
27     criterion,
28     optimizer,
29     scheduler,
30     epochs
31 ):
32     train_losses = []
33     val_losses = []
34
35     for epoch in range(epochs):
36         batch_train_losses = []
37
38         model.train()
39         for idx, inputs in enumerate(train_loader):
40             images = inputs['image']
41             questions = inputs['question']
42             labels = inputs['label']
43
44             optimizer.zero_grad()
45             outputs = model(images, questions)
46             loss = criterion(outputs, labels)
47             loss.backward()
48             optimizer.step()
49
50             batch_train_losses.append(loss.item())
51
52         train_loss = sum(batch_train_losses) / len(batch_train_losses)
53         train_losses.append(train_loss)
54
55         val_loss, val_acc = evaluate(
56             model, val_loader,
57             criterion
58         )
59         val_losses.append(val_loss)
60
61         print(f'EPOCH {epoch + 1}: \tTrain loss: {train_loss:.4f} \tVal loss: {
62             val_loss:.4f} \tVal Acc: {val_acc}')
63
64         scheduler.step()
65
66     return train_losses, val_losses

```

8. Khai báo hàm loss, optimizer và scheduler:

```

1 lr = 1e-2
2 epochs = 50
3 scheduler_step_size = epochs * 0.6
4 criterion = nn.CrossEntropyLoss()
5
6 optimizer = torch.optim.Adam(
7     model.parameters(),
8     lr=lr
9 )
10 scheduler = torch.optim.lr_scheduler.StepLR(

```

```

11     optimizer,
12     step_size=scheduler_step_size,
13     gamma=0.1
14 )
15

```

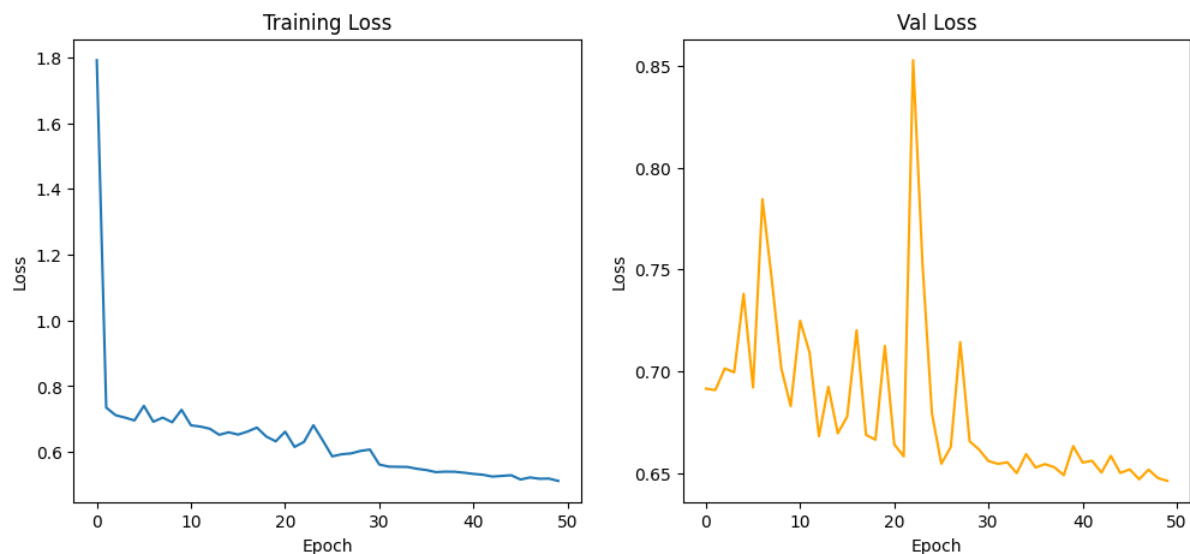
9. **Training:** Với tất cả các thành phần có được, ta tiến hành huấn luyện mô hình:

```

1 train_losses, val_losses = fit(
2     model,
3     train_loader,
4     val_loader,
5     criterion,
6     optimizer,
7     scheduler,
8     epochs
9 )
10

```

Sau khi quá trình huấn luyện hoàn tất, ta có thể trực quan hóa kết quả loss của mô hình trên hai tập train và val qua từng epoch:



Hình 5: Kết quả huấn luyện của mô hình với ViT+RoBERTa

10. **Evaluation:** Cuối cùng, ta đánh giá mô hình đã huấn luyện được trên hai tập val và test:

```

1 val_loss, val_acc = evaluate(
2     model,
3     val_loader,
4     criterion
5 )
6 test_loss, test_acc = evaluate(
7     model,
8     test_loader,
9     criterion
10 )
11
12 print('Evaluation on val/test dataset')
13 print('Val accuracy: ', val_acc)
14 print('Test accuracy: ', test_acc)
15

```

Phần III: Câu hỏi trắc nghiệm

1. Mục tiêu của bài toán Visual Question Answering là gì?
 - (a) Để tăng kích cỡ của tấm ảnh.
 - (b) Để trả lời câu hỏi dựa trên hình ảnh.
 - (c) Để tạo sinh ra hình ảnh dựa trên câu mô tả.
 - (d) Để tạo ra mô hình 3D từ hình ảnh 2D.
2. Trong VQA, mô hình thường nhận đầu vào là gì?
 - (a) Chỉ là hình ảnh
 - (b) Chỉ là câu hỏi văn bản
 - (c) Hình ảnh và câu hỏi văn bản
 - (d) Âm thanh
3. Visual Question Answering là sự kết hợp giữa 2 lĩnh vực nào sau đây?
 - (a) Robotics và Xử lý ngôn ngữ tự nhiên.
 - (b) Thị giác máy tính và học máy.
 - (c) Thị giác máy tính và xử lý ngôn ngữ tự nhiên.
 - (d) Học máy và robotics.
4. Trong VQA, câu trả lời có thể là gì?
 - (a) Chỉ là "có" hoặc "không"
 - (b) Chỉ là một số nguyên
 - (c) Có thể là một câu trả lời văn bản hoặc câu trả lời "có" hoặc "không"
 - (d) Chỉ là một màu sắc
5. Thành phần nào sau đây là một trong những thành phần chính của một mô hình VQA?
 - (a) Nhận diện giọng nói.
 - (b) Trích xuất đặc trưng từ hình ảnh.
 - (c) Dịch câu hỏi ngôn ngữ tự nhiên.
 - (d) Xử lý âm thanh.
6. Một số thách thức trong VQA bao gồm:
 - (a) Hiểu biết ngôn ngữ tự nhiên
 - (b) Nhận dạng đối tượng trong hình ảnh
 - (c) Không có thách thức nào
 - (d) Cả a) và b)
7. Một hệ thống VQA thường xử lý câu hỏi ngôn ngữ tự nhiên như thế nào?
 - (a) Bằng cách chuyển đổi câu hỏi đó thành hình ảnh.
 - (b) Bằng cách áp dụng sentiment analysis cho tấm ảnh đó.
 - (c) Dịch câu hỏi đó qua ngôn ngữ khác.

(d) Bằng cách trích xuất đặc trưng câu hỏi đó bằng các kĩ thuật NLP.

8. Đoạn code sau đây dùng để làm gì?

```
1 eng = spacy.load("en_core_web_sm")
2
3 def get_tokens(data_iter):
4     for sample in data_iter:
5         question = sample['question']
6         yield [token.text for token in eng.tokenizer(question)]
7
```

- (a) Để load mô hình tiếng anh Spacy.
- (b) Tạo ra một danh sách những từ từ những token.
- (c) Thêm ký tự đặc biệt vào các câu hỏi.
- (d) Chia văn bản tiếng Anh thành những token.

9. Mục đích của đoạn code sau là gì?

```
1 vocab = build_vocab_from_iterator(
2     get_tokens(train_data),
3     min_freq=2,
4     specials= ['<pad>', '<sos>', '<eos>', '<unk>'],
5     special_first=True
6 )
7 vocab.set_default_index(vocab['<unk>'])
8
```

- (a) Xây dựng từ vựng từ tập train, bao gồm các token đặc biệt và token mặc định <unk> cho các từ chưa biết.
- (b) Dịch tập data huấn luyện sang một ngôn ngữ khác bằng cách sử dụng các token đặc biệt.
- (c) Sắp xếp các từ trong tập huấn luyện dựa trên tần suất xuất hiện của các từ.
- (d) Đào tạo một mô hình mới để hiểu các mẫu ngôn ngữ trong tập data huấn luyện.

10. Mô hình học máy nào thường được sử dụng cho bài toán VQA?

- (a) Mô hình dựa trên luật có sẵn.
- (b) Các mô hình CNN.
- (c) Mô hình linear regression.
- (d) Mô hình Decision Tree.

11. Thách thức trong bài toán VQA là gì?

- (a) Tăng độ phân giải của hình ảnh
- (b) Sự mơ hồ trong câu hỏi ngôn ngữ tự nhiên
- (c) Xử lý âm thanh.
- (d) Tối ưu hóa phần cứng.

12. Output của mô hình VQA là:

- (a) Bản báo cáo chi tiết.
- (b) Con số cụ thể.
- (c) Câu trả lời cho câu hỏi về hình ảnh.

- (d) Một đồ thị.
13. Hệ thống VQA được ứng dụng trong việc:
- (a) Chơi các video game.
 - (b) Hỗ trợ người dùng khiếm thị trong việc hiểu môi trường xung quanh.
 - (c) Lọc email rác.
 - (d) Xử lý âm thanh.
14. Dòng code sau đây có tác dụng gì?
- ```
1 image_encoder = timm.create_model(image_model, pretrained=True, num_classes=
 hidden_dim))
2
```
- (a) Tải pretrained hình ảnh từ thư viện timm.
  - (b) Chuyển đổi dữ liệu hình ảnh màu sang hình ảnh xám.
  - (c) Tăng độ phân giải của hình ảnh.
  - (d) Nén hình ảnh để xử lý nhanh hơn.
15. Output của image\_encoder trong đoạn code sau là gì?
- ```
1 image_encoder = timm.create_model(image_model, pretrained=True, num_classes=  
    hidden_dim))  
2
```
- (a) Ảnh được biến đổi từ tập dữ liệu.
 - (b) Một vector embedding với số chiều là hidden_dim
 - (c) Xác suất tầm ảnh đầu vào thuộc về từng class.
 - (d) Chiều của tấm ảnh đầu vào.
16. Điều KHÔNG làm ảnh hưởng đến độ chính xác của một mô hình VQA?
- (a) Tốc độ Internet.
 - (b) Chất lượng của tập dữ liệu đầu vào.
 - (c) Độ nhiễu trong các tấm ảnh đầu vào.
 - (d) Độ phức tạp trong câu hỏi từ ngôn ngữ tự nhiên.
17. Công nghệ nào là cần thiết để giải thích thành phần văn bản trong VQA?
- (a) Blockchain
 - (b) Xử lý Ngôn Ngữ Tự Nhiên (NLP)
 - (c) Tính Toán Lượng Tử
 - (d) Mật mã học

- Hết -