

## ✓ MOUNT GOOGLE DRIVE

```

1 # Mount Google Drive
2 from google.colab import drive
3 import os
4
5 # Mount Google Drive
6 drive.mount('/content/drive')

```

➞ Mounted at /content/drive

## ✓ 1. EDA and Data Processing

```

1 import os
2 import json
3 import pickle
4 import torch
5 from PIL import Image
6 from torch.utils.data import Dataset, DataLoader
7 from transformers import BertTokenizer, BertModel
8 from torchvision import transforms
9 import torch.nn as nn
10 import matplotlib.pyplot as plt
11
12 # ✓ Paths
13 base_path = '/content/drive/MyDrive/VQA_Project'
14 image_folder = os.path.join(base_path, 'Images/AbstractScenes_v1')
15
16 # ✓ Load questions & annotations
17 with open(os.path.join(base_path, 'Questions/MultipleChoice_Questions_train_abstract_v1.json')) as f:
18     train_questions = json.load(f)['questions']
19 with open(os.path.join(base_path, 'Questions/MultipleChoice_Questions_val_abstract_v1.json')) as f:
20     val_questions = json.load(f)['questions']
21 with open(os.path.join(base_path, 'Annotations/Abstract_Scenes_Annotations_Train.json')) as f:
22     train_annotations = json.load(f)['annotations']
23 with open(os.path.join(base_path, 'Annotations/Abstract_Scenes_Annotations_Val.json')) as f:
24     val_annotations = json.load(f)['annotations']
25
26 # ✓ Prepare Data
27 train_ans_map = {ann['question_id']: ann['multiple_choice_answer'] for ann in train_annotations}
28 val_ans_map = {ann['question_id']: ann['multiple_choice_answer'] for ann in val_annotations}
29
30 def prepare_data(questions, ans_map, split='train'):
31     data = []
32     for q in questions:

```

```

33     q_id = q['question_id']
34     image_id = q['image_id']
35     img_filename = f"abstract_v002_{split}2015_{image_id:012d}.png"
36     if ans_map.get(q_id) in q['multiple_choices']:
37         data.append({
38             'image': img_filename,
39             'question': q['question'],
40             'choices': q['multiple_choices'],
41             'answer': ans_map[q_id]
42         })
43     return data
44
45 train_data = prepare_data(train_questions, train_ans_map, split='train')
46 val_data = prepare_data(val_questions, val_ans_map, split='val')
47
48 # ✓ Image transform
49 image_transform = transforms.Compose([
50     transforms.Resize((224, 224)),
51     transforms.ToTensor()
52 ])
53
54 # ✓ Tokenizer
55 tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
56
57 # ✓ Dataset class with safety patch
58 class VQADataset(Dataset):
59     def __init__(self, data, image_folder, tokenizer, transform=None, max_length=32):
60         self.data = data
61         self.image_folder = image_folder
62         self.tokenizer = tokenizer
63         self.transform = transform
64         self.max_length = max_length
65
66     def __len__(self):
67         return len(self.data)
68
69     def __getitem__(self, idx):
70         item = self.data[idx]
71         image_path = os.path.join(self.image_folder, item["image"])
72         try:
73             image = Image.open(image_path).convert("RGB")
74         except:
75             image = Image.new("RGB", (224, 224), color=(0, 0, 0)) # Black placeholder
76
77         if self.transform:
78             image = self.transform(image)
79
80         encoded = self.tokenizer(
81             [f"{item['question']} {choice}" for choice in item["choices"]],
82             padding="max_length", truncation=True, max_length=self.max_length, return_t
83 )

```

```

84
85     label = item["choices"].index(item["answer"])
86
87     return {
88         "image": image,
89         "input_ids": encoded["input_ids"],
90         "attention_mask": encoded["attention_mask"],
91         "label": torch.tensor(label),
92         "question": item["question"], # ✓ Include question
93         "choices": item["choices"],   # ✓ Include choices
94         "answer": item["answer"]      # ✓ Include correct answer
95     }
96
97 train_dataset = VQADataset(train_data, image_folder, tokenizer, transform=image_transfc
98 val_dataset = VQADataset(val_data, image_folder, tokenizer, transform=image_transform)
99 train_loader = DataLoader(train_dataset, batch_size=32, shuffle=True, num_workers=0)
100 val_loader = DataLoader(val_dataset, batch_size=32, shuffle=False, num_workers=0)
101
102 /usr/local/lib/python3.11/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)
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 model
warnings.warn(
tokenizer_config.json: 100% 48.0/48.0 [00:00<00:00, 5.70kB/s]
vocab.txt: 100% 232k/232k [00:00<00:00, 4.59MB/s]
tokenizer.json: 100% 466k/466k [00:00<00:00, 23.1MB/s]
config.json: 100% 570/570 [00:00<00:00, 72.7kB/s]

```

## Check Last Saved Model & Accuracy

```

1 import os
2 import pickle
3 import datetime
4
5 # ✓ Set paths
6 base_path = "/content/drive/MyDrive/VQA_Project"
7 model_path = os.path.join(base_path, "vqa_finetuned_best.pt")
8 history_path = os.path.join(base_path, "train_history.pkl")
9
10 # ✓ Check last modified time
11 if os.path.exists(model_path):
12     last_modified = os.path.getmtime(model_path)
13     last_save_time = datetime.datetime.fromtimestamp(last_modified).strftime('%Y-%m-%d %H:%M:%S')
14     print(f"📁 Model Last Saved: {last_save_time}")
15 else:
16     print("❌ Model file not found.")

```

```

17
18 # ✅ Load last saved validation accuracy
19 if os.path.exists(history_path):
20     with open(history_path, "rb") as f:
21         history = pickle.load(f)
22
23     # Extract last epoch data
24     if "val accuracies" in history and history["val accuracies"]:
25         last_epoch = len(history["val accuracies"])
26         last_val_acc = history["val accuracies"][-1] * 100 # Convert to percentage
27         last_train_acc = history["train accuracies"][-1] * 100 # Convert to percentage
28
29         print(f"📁 Last Saved Epoch: {last_epoch}")
30         print(f"✅ Last Train Accuracy: {last_train_acc:.2f}%")
31         print(f"✅ Last Validation Accuracy: {last_val_acc:.2f}%")
32     else:
33         print("⚠️ No accuracy data found in history file.")
34 else:
35     print("❌ Training history file not found.")
36
📁 Model Last Saved: 2025-03-24 13:53:03
📁 Last Saved Epoch: 8
✅ Last Train Accuracy: 75.67%
✅ Last Validation Accuracy: 63.72%

```

## ✓ Load the Model

```

1 import torch
2 import torch.nn as nn
3 from transformers import BertModel
4 from torchvision.models import resnet18
5
6 class VQAModel(nn.Module):
7     def __init__(self, hidden_dim=512, num_choices=4): # Ensure num_choices is set corr
8         super(VQAModel, self).__init__()
9
10        # ✅ Image encoder (ResNet18 without final classifier)
11        self.cnn = resnet18(pretrained=True)
12        self.cnn.fc = nn.Identity() # Remove classification layer
13        self.cnn_out_dim = 512 # ResNet18 output size
14
15        # ✅ Text encoder (BERT, frozen)
16        self.bert = BertModel.from_pretrained('bert-base-uncased')
17        for param in self.bert.parameters():
18            param.requires_grad = False # Freeze BERT during training
19
20        self.bert_out_dim = self.bert.config.hidden_size # 768
21
22        # ✅ Fusion layer (Image + Text features)

```

```

23         self.fusion = nn.Sequential(
24             nn.Linear(self.cnn_out_dim + self.bert_out_dim, hidden_dim),
25             nn.ReLU(),
26             nn.Dropout(0.3),
27             nn.Linear(hidden_dim, 1) # One score per choice
28         )
29
30     def forward(self, images, input_ids, attention_mask):
31         batch_size, num_choices, seq_len = input_ids.shape
32
33         # ✅ Encode images
34         img_feat = self.cnn(images) # [batch, 512]
35         img_feat = img_feat.unsqueeze(1).repeat(1, num_choices, 1) # [batch, num_choices, 512]
36
37         # ✅ Encode text
38         input_ids = input_ids.view(-1, seq_len) # Flatten for BERT batch processing
39         attention_mask = attention_mask.view(-1, seq_len)
40         text_outputs = self.bert(input_ids=input_ids, attention_mask=attention_mask)
41         text_feat = text_outputs.pooler_output # [batch*num_choices, 768]
42         text_feat = text_feat.view(batch_size, num_choices, -1) # [batch, num_choices, 768]
43
44         # ✅ Fuse image + text
45         fused = torch.cat([img_feat, text_feat], dim=-1) # [batch, num_choices, 1280]
46         scores = self.fusion(fused).squeeze(-1) # [batch, num_choices]
47
48         return scores # Logits
49

```

## 2 Load Model Checkpoint

```

1 import torch
2 import os
3 import torchvision.models as models
4 from torchvision.models import ResNet18_Weights
5
6 resnet18(weights=ResNet18_Weights.DEFAULT) # For most up-to-date weights
7 backbone = resnet18(weights=ResNet18_Weights.DEFAULT)
8 backbone.fc = nn.Identity()
9
10 # ✅ Define device
11 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
12
13 # ✅ Initialize model & move to device
14 model = VQAModel().to(device)
15
16 # ✅ Load the checkpoint with correct architecture
17 checkpoint_path = os.path.join(base_path, "vqa_finetuned_best.pt")
18 state_dict = torch.load(checkpoint_path, map_location=device, weights_only=True)
19
20 # ✅ Load state dict with `strict=False` to avoid minor mismatches

```



```

6 # ✅ Configure Matplotlib to use fonts available on your system
7 from matplotlib import rcParams
8 rcParams['font.family'] = ['DejaVu Sans', 'Liberation Sans', 'STIXGeneral'] # Based on
9 rcParams['axes.unicode_minus'] = False # Ensure proper rendering of minus signs and Uni
10
11 # ✅ Define function to randomly pick images from the dataset for sample predictions
12 def show_random_predictions(model, dataset, num_samples=10):
13     model.eval()
14     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
15
16     fig, axes = plt.subplots(num_samples, 1, figsize=(8, 6 * num_samples)) # Increased
17     if num_samples == 1:
18         axes = [axes] # Ensure iterable format
19
20     # 🎲 Randomly select samples from the dataset
21     sample_indices = random.sample(range(len(dataset)), num_samples)
22
23     with torch.no_grad():
24         for i, idx in enumerate(sample_indices):
25             batch = dataset[idx] # Directly fetch a sample from dataset
26             image = batch['image'].unsqueeze(0).to(device) # Add batch dimension
27             input_ids = batch['input_ids'].unsqueeze(0).to(device)
28             attention_mask = batch['attention_mask'].unsqueeze(0).to(device)
29             label = batch['label'].item()
30
31             # 🔮 Make prediction
32             logits = model(image, input_ids, attention_mask)
33             pred_idx = torch.argmax(logits, dim=1).item()
34
35             # 🖼️ Get question & choices
36             img = image[0].cpu().permute(1, 2, 0).numpy()
37             question = batch['question']
38             choices = batch['choices']
39
40             # 🔄 Ensure index stays within bounds
41             predicted_index = min(pred_idx, len(choices) - 1)
42             correct_index = min(label, len(choices) - 1)
43
44             predicted_choice = choices[predicted_index]
45             correct_choice = choices[correct_index]
46
47             # ✅ Text-based alternatives for Emojis
48             correct_emoji = "✅" # Unicode check mark
49             predicted_emoji = "🔮" # Unicode crystal ball
50
51             # 🎨 Plot
52             axes[i].imshow(img)
53             axes[i].axis("off")
54             axes[i].set_title(
55                 f"Q: {question}\n{correct_emoji} Correct: {correct_choice} | {predicted_emoji} Predicted: {predicted_choice}"
56                 color="green" if predicted_choice == correct_choice else "red"

```

```
57         )
58
59     plt.tight_layout()
60     plt.show()
61
62 # 🔥 Run with randomized selection from the dataset
```

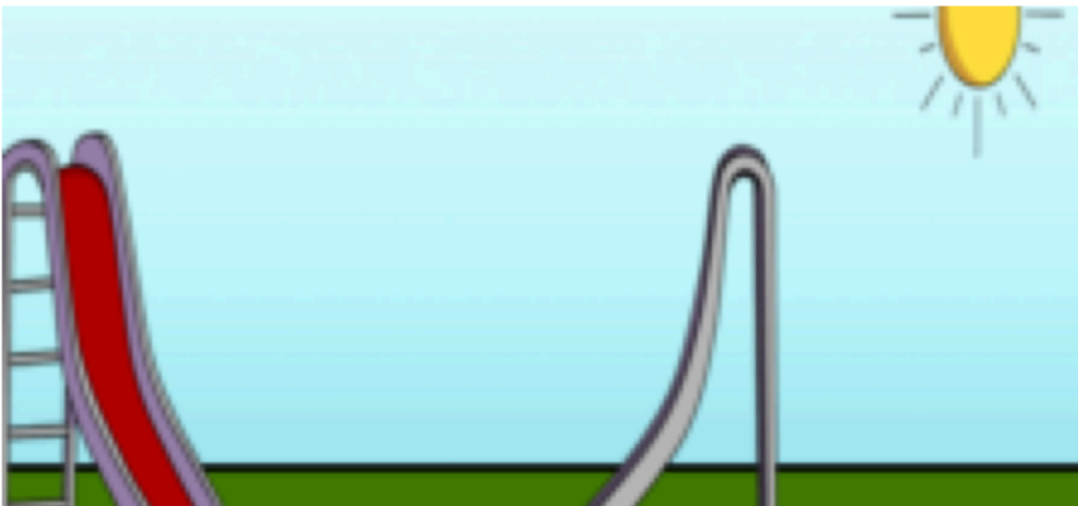


```
>>> <ipython-input-8-e9719dd27217>:59: UserWarning: Glyph 128302 (\N{CRYSTAL BALL}) missing  
plt.tight_layout()  
/usr/local/lib/python3.11/dist-packages/IPython/core/pylabtools.py:151: UserWarning: Gly  
fig.canvas.print_figure(bytes_io, **kw)
```

Q: Is he angry at the cat?  
✓ Correct: yes | ☐ Predicted: no



Q: Is it sunny?  
✓ Correct: yes | ☐ Predicted: yes





Q: Where is the game control at?  
✓ Correct: floor | ☐ Predicted: floor



Q: What is being kicked?  
✓ Correct: soccer ball | ☐ Predicted: football



Q: Is the woman going to eat alone?

✓ Correct: yes | ☐ Predicted: yes



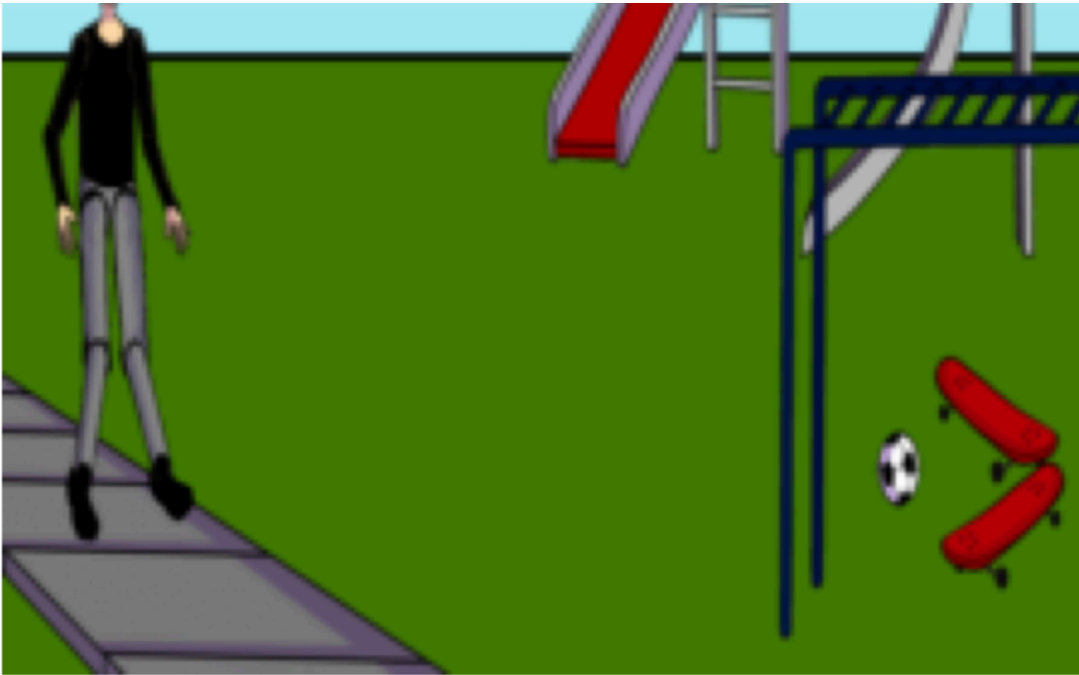


Q: Is there a TV?  
✓ Correct: no | ☐ Predicted: no



Q: Is that a basketball on the ground?  
✓ Correct: no | ☐ Predicted: no





Q: Are they women?  
✓ Correct: yes | ☐ Predicted: no



Q: Does the dog belong to the woman?  
✓ Correct: yes | ☐ Predicted: yes

✓ Correct: yes | ☐ Predicted: yes



Q: Are the 2 people dancing?  
✓ Correct: yes | ☐ Predicted: yes





## ✅ Step 2: Generate a Classification Report

```

1 from sklearn.metrics import classification_report
2
3 # ✅ Function to evaluate and generate classification report
4 def generate_classification_report(model, dataloader):
5     model.eval()
6     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
7
8     all_preds, all_labels = [], []
9
10    with torch.no_grad():
11        for batch in dataloader:
12            images = batch['image'].to(device)
13            input_ids = batch['input_ids'].to(device)
14            attention_mask = batch['attention_mask'].to(device)
15            labels = batch['label'].cpu().numpy() # Convert to NumPy
16
17            logits = model(images, input_ids, attention_mask)
18            preds = torch.argmax(logits, dim=1).cpu().numpy()
19
20            all_preds.extend(preds)
21            all_labels.extend(labels)
22
23    # 📄 Print classification report
24    print("📄 Classification Report:")
25    print(classification_report(all_labels, all_preds))
26
27 # 🔥 Run classification report
28 generate_classification_report(model, val_loader)
29

```



📄 Classification Report:

	precision	recall	f1-score	support
0	0.64	0.63	0.63	1715
1	0.65	0.66	0.65	1671
2	0.62	0.64	0.63	1652
3	0.65	0.64	0.65	1681
4	0.63	0.63	0.63	1623
5	0.63	0.64	0.63	1657
6	0.65	0.65	0.65	1694
7	0.64	0.64	0.64	1687
8	0.64	0.64	0.64	1685
9	0.65	0.63	0.64	1641
10	0.64	0.64	0.64	1668
11	0.64	0.64	0.64	1673
12	0.65	0.64	0.64	1669
13	0.64	0.64	0.64	1656
14	0.65	0.65	0.65	1646
15	0.64	0.63	0.63	1641
16	0.64	0.63	0.63	1640



	17	0.62	0.63	0.63	1701
accuracy				0.64	30000
macro avg		0.64	0.64	0.64	30000
weighted avg		0.64	0.64	0.64	30000

1 Start coding or [generate](#) with AI.

1 Start coding or [generate](#) with AI.

## ✅ Step 1: Extract True Labels (y\_true) & Predictions (y\_pred)

```

1 import torch
2 from torch.utils.data import DataLoader
3
4 val_dataloader = DataLoader(val_dataset, batch_size=32, shuffle=False) # Adjust batch_s
5
6 y_true = []
7 y_pred = []
8
9 # Ensure model is in evaluation mode
10 model.eval()
11 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
12
13 # Loop through validation data
14 with torch.no_grad():
15     for batch in val_dataloader: # Replace with your actual DataLoader
16         images = batch['image'].to(device)
17         input_ids = batch['input_ids'].to(device)
18         attention_mask = batch['attention_mask'].to(device)
19         labels = batch['label'].cpu().numpy() # Extract true labels
20
21         # Get predictions
22         logits = model(images, input_ids, attention_mask)
23         predictions = torch.argmax(logits, dim=1).cpu().numpy() # Convert to numpy
24
25         # Store results
26         y_true.extend(labels)
27         y_pred.extend(predictions)
28

```

```

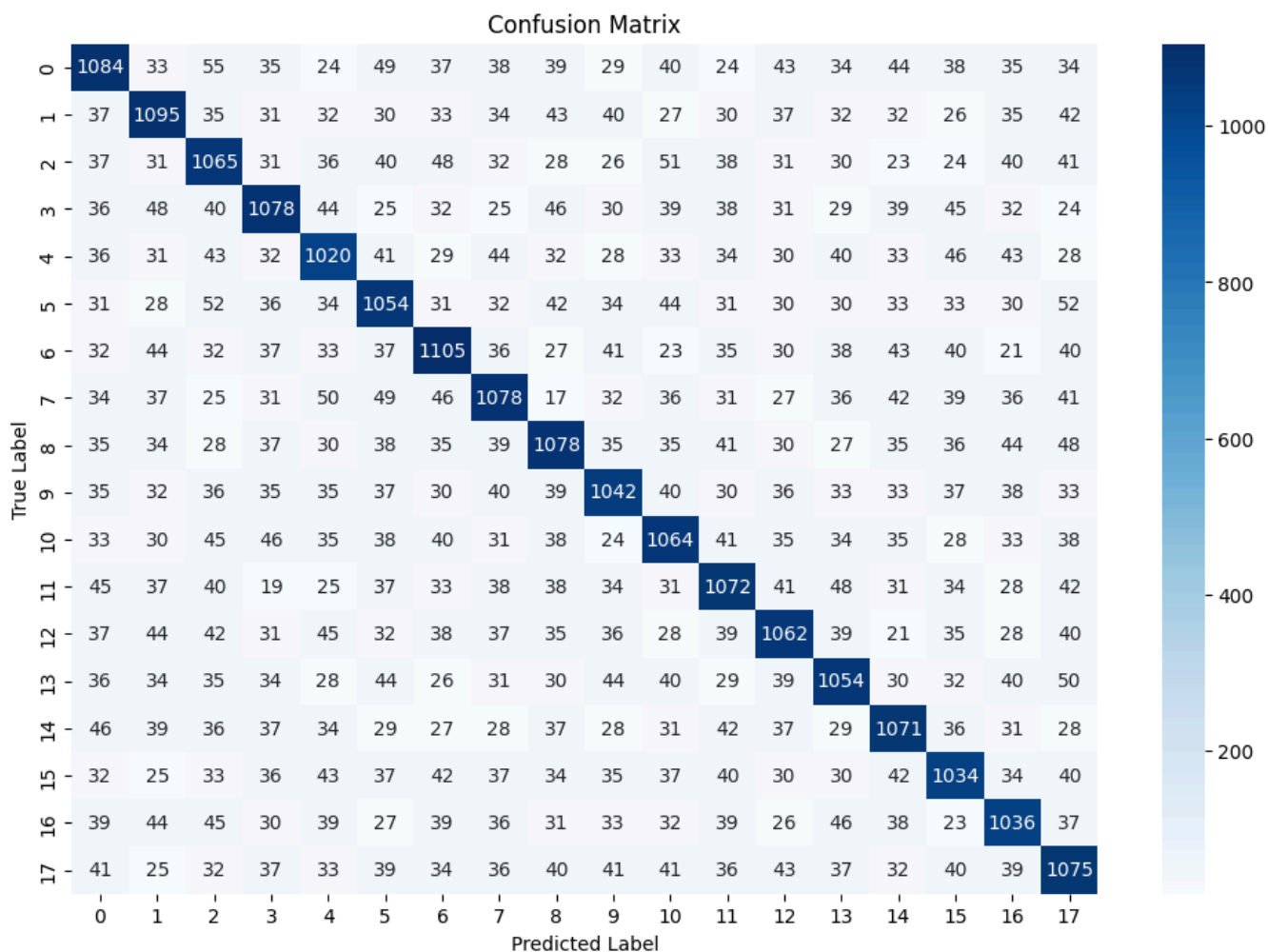
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 from sklearn.metrics import confusion_matrix
4
5 # Compute confusion matrix
6 cm = confusion_matrix(y_true, y_pred)

```


```

7
8 # Get class labels (if available)
9 class_labels = sorted(set(y_true)) # Extract unique class labels
10
11 # Plot confusion matrix
12 plt.figure(figsize=(12, 8))
13 sns.heatmap(cm, annot=True, fmt="d", cmap="Blues", xticklabels=class_labels, yticklabels=
14 plt.xlabel("Predicted Label")
15 plt.ylabel("True Label")
16 plt.title("Confusion Matrix")
17 plt.show()
18

```



```
1 from collections import Counter
2
3 # Initialize counters
4 question_type_counter = Counter()
5 answer_counter = Counter()
6
7 # Single iteration over the dataset
8 for sample in train_dataset:
9     question_type_counter.update([sample['question'][:10].lower()]) # Count question ty
10     answer_counter.update([sample['answer']]) # Count answers
11
12 # Display results
13 print("Top question types:", question_type_counter.most_common(10))
14 print("Answer distribution:", answer_counter.most_common(10))
```

```
1 import pickle
2 import matplotlib.pyplot as plt
3
4 #  Load training history
5 with open("/content/drive/MyDrive/VQA_Project/train_history.pkl", "rb") as f:
6     history = pickle.load(f)
7
8 train_losses = history["train_losses"]
9 val_losses = history["val_losses"]
10 train_accuracies = history["train_accuracies"]
11 val_accuracies = history["val_accuracies"]
```