```python
1  import torch
2  import torch.nn as nn
3  import torch.optim as optim
4  from torch.utils.data import Dataset, DataLoader
5  from collections import defaultdict
6  import random
7
8
9  def read_tsv(path):
10     data = []
11     with open(path, encoding='utf-8') as f:
12         for line in f:
13             dev, lat, freq = line.strip().split('\t')
14             data.extend([(lat, dev)] * int(freq))
15     return data
16
17 def build_vocab(sequences):
18     vocab = {'<pad>': 0, '<sos>': 1, '<eos>': 2}
19     for seq in sequences:
20         for char in seq:
21             if char not in vocab:
22                 vocab[char] = len(vocab)
23     return vocab
24
25
26 class TransliterationDataset(Dataset):
27     def __init__(self, data, input_vocab, target_vocab):
28         self.data = data
29         self.input_vocab = input_vocab
30         self.target_vocab = target_vocab
31
32     def __len__(self):
33         return len(self.data)
34
35     def encode_seq(self, seq, vocab, add_sos_eos=False):
36         ids = [vocab[c] for c in seq]
37         if add_sos_eos:
38             ids = [vocab['<sos>']] + ids + [vocab['<eos>']]
39         return torch.tensor(ids, dtype=torch.long)
40
41     def __getitem__(self, idx):
42         latin, dev = self.data[idx]
43         return self.encode_seq(latin, self.input_vocab), self.encode_seq(dev, self.target_vocab, True)
44
45 def collate_fn(batch):
46     srcs, trgs = zip(*batch)
47     srcs_padded = nn.utils.rnn.pad_sequence(srcs, batch_first=True, padding_value=0)
48     trgs_padded = nn.utils.rnn.pad_sequence(trgs, batch_first=True, padding_value=0)
49     return srcs_padded, trgs_padded
50
51
52 class Encoder(nn.Module):
53     def __init__(self, input_dim, emb_dim, hid_dim, n_layers, rnn_type='gru'):
54         super().__init__()
55         self.embedding = nn.Embedding(input_dim, emb_dim)
56         rnn_cls = {'rnn': nn.RNN, 'lstm': nn.LSTM, 'gru': nn.GRU}[rnn_type]
57         self.rnn = rnn_cls(emb_dim, hid_dim, n_layers, batch_first=True)
58         self.rnn_type = rnn_type
59
60     def forward(self, src):
61         embedded = self.embedding(src)
62         outputs, hidden = self.rnn(embedded)
63         return hidden
64
65 class Decoder(nn.Module):
66     def __init__(self, output_dim, emb_dim, hid_dim, n_layers, rnn_type='gru'):
67         super().__init__()
68         self.embedding = nn.Embedding(output_dim, emb_dim)
69         rnn_cls = {'rnn': nn.RNN, 'lstm': nn.LSTM, 'gru': nn.GRU}[rnn_type]
70         self.rnn = rnn_cls(emb_dim, hid_dim, n_layers, batch_first=True)
71         self.fc_out = nn.Linear(hid_dim, output_dim)
72
73     def forward(self, input, hidden):
74         input = input.unsqueeze(1)
75         embedded = self.embedding(input)
76         output, hidden = self.rnn(embedded, hidden)
77         prediction = self.fc_out(output.squeeze(1))
78         return prediction, hidden
79
80 class Seq2Seq(nn.Module):
81     def __init__(self, encoder, decoder, device):
82         super().__init__()
83         self.encoder = encoder
84         self.decoder = decoder
```

```
 85            self.device = device
 86
 87     def forward(self, src, trg, teacher_forcing_ratio=0.5):
 88            batch_size, trg_len = trg.shape
 89            output_dim = self.decoder.fc_out.out_features
 90            outputs = torch.zeros(batch_size, trg_len, output_dim).to(self.device)
 91
 92            hidden = self.encoder(src)
 93            input = trg[:, 0]
 94
 95            for t in range(1, trg_len):
 96                output, hidden = self.decoder(input, hidden)
 97                outputs[:, t] = output
 98                teacher_force = random.random() < teacher_forcing_ratio
 99                top1 = output.argmax(1)
100                input = trg[:, t] if teacher_force else top1
101
102            return outputs
103
104
105 def train(model, data_loader, optimizer, criterion, clip=1):
106     model.train()
107     epoch_loss = 0
108     for src, trg in data_loader:
109         src, trg = src.to(model.device), trg.to(model.device)
110         optimizer.zero_grad()
111         output = model(src, trg)
112         output_dim = output.shape[-1]
113         output = output[:, 1:].reshape(-1, output_dim)
114         trg = trg[:, 1:].reshape(-1)
115         loss = criterion(output, trg)
116         loss.backward()
117         torch.nn.utils.clip_grad_norm_(model.parameters(), clip)
118         optimizer.step()
119         epoch_loss += loss.item()
120     return epoch_loss / len(data_loader)
121
122 def accuracy(model, data_loader):
123     model.eval()
124     correct, total = 0, 0
125     with torch.no_grad():
126         for src, trg in data_loader:
127             src, trg = src.to(model.device), trg.to(model.device)
128             output = model(src, trg, 0)
129             preds = output.argmax(-1)
130             for pred, true in zip(preds, trg):
131                 if torch.equal(pred[1:], true[1:]):
132                     correct += 1
133                 total += 1
134     return correct / total
135
136 def predict(model, src_seq, input_vocab, output_vocab, max_len=30):
137     model.eval()
138     inv_vocab = {v: k for k, v in output_vocab.items()}
139     src_tensor = torch.tensor([input_vocab[c] for c in src_seq], dtype=torch.long).unsqueeze(0).to(model.device)
140     hidden = model.encoder(src_tensor)
141     input = torch.tensor([output_vocab['<sos>']], device=model.device)
142     output = []
143     for _ in range(max_len):
144         out, hidden = model.decoder(input, hidden)
145         top1 = out.argmax(1)
146         char = inv_vocab[top1.item()]
147         if char == '<eos>':
148             break
149         output.append(char)
150         input = top1
151     return ''.join(output)
152
153
154 DEVICE = torch.device("cuda" if torch.cuda.is_available() else "cpu")
155
156 train_path = "/content/hi.translit.sampled.train.tsv"
157 val_path = "/content/hi.translit.sampled.dev.tsv"
158 train_data = read_tsv(train_path)
159 val_data = read_tsv(val_path)
160
161 input_vocab = build_vocab([d[0] for d in train_data])
162 target_vocab = build_vocab([d[1] for d in train_data])
163
164 train_dataset = TransliterationDataset(train_data, input_vocab, target_vocab)
165 val_dataset = TransliterationDataset(val_data, input_vocab, target_vocab)
166 train_loader = DataLoader(train_dataset, batch_size=32, shuffle=True, collate_fn=collate_fn)
167 val_loader = DataLoader(val_dataset, batch_size=32, shuffle=False, collate_fn=collate_fn)
168
169 INPUT_DIM = len(input_vocab)
170 OUTPUT_DIM = len(target_vocab)
```

```python
171 EMB_DIM = 64
172 HID_DIM = 128
173 N_LAYERS = 1
174 RNN_TYPE = 'gru'
175
176 encoder = Encoder(INPUT_DIM, EMB_DIM, HID_DIM, N_LAYERS, RNN_TYPE)
177 decoder = Decoder(OUTPUT_DIM, EMB_DIM, HID_DIM, N_LAYERS, RNN_TYPE)
178 model = Seq2Seq(encoder, decoder, DEVICE).to(DEVICE)
179
180 optimizer = optim.Adam(model.parameters(), lr=0.001)
181 criterion = nn.CrossEntropyLoss(ignore_index=target_vocab['<pad>'])
182
183 for epoch in range(10):
184     loss = train(model, train_loader, optimizer, criterion)
185     acc = accuracy(model, val_loader)
186     print(f"Epoch {epoch+1} | Loss: {loss:.4f} | Val Accuracy: {acc:.4f}")
187
188 print("\nSample Predictions:")
189 for i in range(5):
190     src_sample, tgt_sample = val_data[i]
191     pred = predict(model, src_sample, input_vocab, target_vocab)
192     print(f"Input: {src_sample} | Target: {tgt_sample} | Predicted: {pred}")
193
```

```
Epoch 1 | Loss: 1.5929 | Val Accuracy: 0.0069
Epoch 2 | Loss: 0.8821 | Val Accuracy: 0.0103
Epoch 3 | Loss: 0.7274 | Val Accuracy: 0.0162
Epoch 4 | Loss: 0.6378 | Val Accuracy: 0.0167
Epoch 5 | Loss: 0.5878 | Val Accuracy: 0.0204
Epoch 6 | Loss: 0.5420 | Val Accuracy: 0.0238
Epoch 7 | Loss: 0.5066 | Val Accuracy: 0.0199
Epoch 8 | Loss: 0.4810 | Val Accuracy: 0.0216
Epoch 9 | Loss: 0.4545 | Val Accuracy: 0.0219
Epoch 10 | Loss: 0.4318 | Val Accuracy: 0.0232

Sample Predictions:
Input: ankan | Target: अंकन | Predicted: अंकान
Input: ankan | Target: अंकन | Predicted: अंकान
Input: ankan | Target: अंकन | Predicted: अंकान
Input: angkor | Target: अंगकोर | Predicted: अंगकर
Input: angkor | Target: अंगकोर | Predicted: अंगकर
```

```python
1 import pandas as pd
2 import re
3 from datasets import load_dataset
4 from transformers import GPT2Tokenizer, GPT2LMHeadModel, Trainer, TrainingArguments, pipeline
5
6
7 def load_and_clean_lyrics(*csv_paths):
8     dfs = [pd.read_csv(path) for path in csv_paths]
9     lyrics_df = pd.concat(dfs)
10
11     def clean_lyrics(lyric):
12         if pd.isna(lyric):
13             return ""
14         lyric = str(lyric)
15         lyric = re.sub(r'^#+', '', lyric)
16         lyric = lyric.encode('utf-8').decode('utf-8', 'ignore')
17         lyric = re.sub(r'[\u2018\u2019\u201c\u201d]+', "'", lyric)
18         lyric = re.sub(r'[^\x00-\x7F]+', '', lyric)
19         return lyric.strip()
20
21     return lyrics_df['Lyric'].dropna().apply(clean_lyrics).tolist()
22
23 lyrics_texts = load_and_clean_lyrics('/content/EdSheeran.csv', '/content/JustinBieber.csv')
24
25 with open("lyrics_dataset.txt", "w", encoding="utf-8") as f:
26     for lyric in lyrics_texts:
27         f.write(lyric + "\n\n")
28
29 dataset = load_dataset("text", data_files={"train": "lyrics_dataset.txt"})
30
31 tokenizer = GPT2Tokenizer.from_pretrained("gpt2")
32 tokenizer.pad_token = tokenizer.eos_token
33
34 def tokenize_function(example):
35     return tokenizer(example["text"], truncation=True, padding="max_length", max_length=512)
36
37 tokenized_dataset = dataset.map(tokenize_function, batched=True, remove_columns=["text"])
38 tokenized_dataset = tokenized_dataset.map(lambda examples: {"labels": examples["input_ids"]}, batched=True)
39
40 model = GPT2LMHeadModel.from_pretrained("gpt2")
41
42 training_args = TrainingArguments(
43     output_dir="./gpt2-lyrics",
44     per_device_train_batch_size=2,
```

```python
45      num_train_epochs=3,
46      logging_steps=100,
47      save_steps=500,
48      save_total_limit=1,
49      prediction_loss_only=True,
50      report_to="none",
51      fp16=False
52  )
53
54  trainer = Trainer(
55      model=model,
56      args=training_args,
57      train_dataset=tokenized_dataset["train"],
58      tokenizer=tokenizer
59  )
60
61  trainer.train()
62
63
64
65  generator = pipeline("text-generation", model=model, tokenizer=tokenizer)
66  output = generator("I remember those nights when", max_length=100, num_return_sequences=1)[0]["generated_text"]
67  print(output)
68
```

Generating train split:  ▮  1282/0 [00:00<00:00, 34168.72 examples/s]
/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), set it as secret in yo
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, 3.00kB/s]
vocab.json: 100% ▮▮▮▮▮▮▮▮▮▮▮▮▮▮  1.04M/1.04M [00:00<00:00, 5.18MB/s]
merges.txt: 100% ▮▮▮▮▮▮▮▮▮▮▮▮▮▮  456k/456k [00:00<00:00, 3.18MB/s]
tokenizer.json: 100% ▮▮▮▮▮▮▮▮▮▮▮▮▮  1.36M/1.36M [00:00<00:00, 9.08MB/s]
config.json: 100% ▮▮▮▮▮▮▮▮▮▮▮▮▮  665/665 [00:00<00:00, 80.3kB/s]
Map: 100% ▮▮▮▮▮▮▮▮▮▮▮▮▮  1282/1282 [00:01<00:00, 714.99 examples/s]
Map: 100% ▮▮▮▮▮▮▮▮▮▮▮▮▮  1282/1282 [00:00<00:00, 2836.24 examples/s]
Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular HTTP download. For better performance,
WARNING:huggingface_hub.file_download:Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular
model.safetensors: 100% ▮▮▮▮▮▮▮▮▮▮▮▮▮  548M/548M [00:01<00:00, 310MB/s]
generation_config.json: 100% ▮▮▮▮▮▮▮▮▮▮▮▮▮  124/124 [00:00<00:00, 8.76kB/s]
<ipython-input-5-c9ee5be06c98>:70: FutureWarning: `tokenizer` is deprecated and will be removed in version 5.0.0 for `Trainer.__init__`. Use `p
  trainer = Trainer(