```
1 import pandas as pd
2
3 df=pd.read_csv('/content/CEAS_08[1].csv')
4 df.head()
```

| | sender | receiver | date | subject |
|---|---|---|---|---|
| 0 | Young Esposito <Young@iworld.de> | user4@gvc.ceas-challenge.cc | Tue, 05 Aug 2008 16:31:02 -0700 | Never agree to be a loser |
| 1 | Mok <ipline's1983@icable.ph> | user2.2@gvc.ceas-challenge.cc | Tue, 05 Aug 2008 18:31:03 -0500 | Befriend Jenna Jameson |
| 2 | Daily Top 10 <Karmandeep-opengevl@universalnet... | user2.9@gvc.ceas-challenge.cc | Tue, 05 Aug 2008 20:28:00 -1200 | CNN.com Daily Top 10   >+=+ |
| 3 | Michael Parker <ivqrnai@pobox.com> | SpamAssassin Dev <xrh@spamassassin.apache.org> | Tue, 05 Aug 2008 17:31:20 -0600 | Re: svn commit: r619753 - in /spamassassin/tru... |
| 4 | Gretchen Suggs <externalsep1@loanofficertool.com> | user2.2@gvc.ceas-challenge.cc | Tue, 05 Aug 2008 19:31:21 -0400 | SpecialPricesPharmMoreinfo |

```
1 import pandas as pd
2 df=pd.read_csv('/content/CEAS_08[1].csv')
3 df.size
4
```
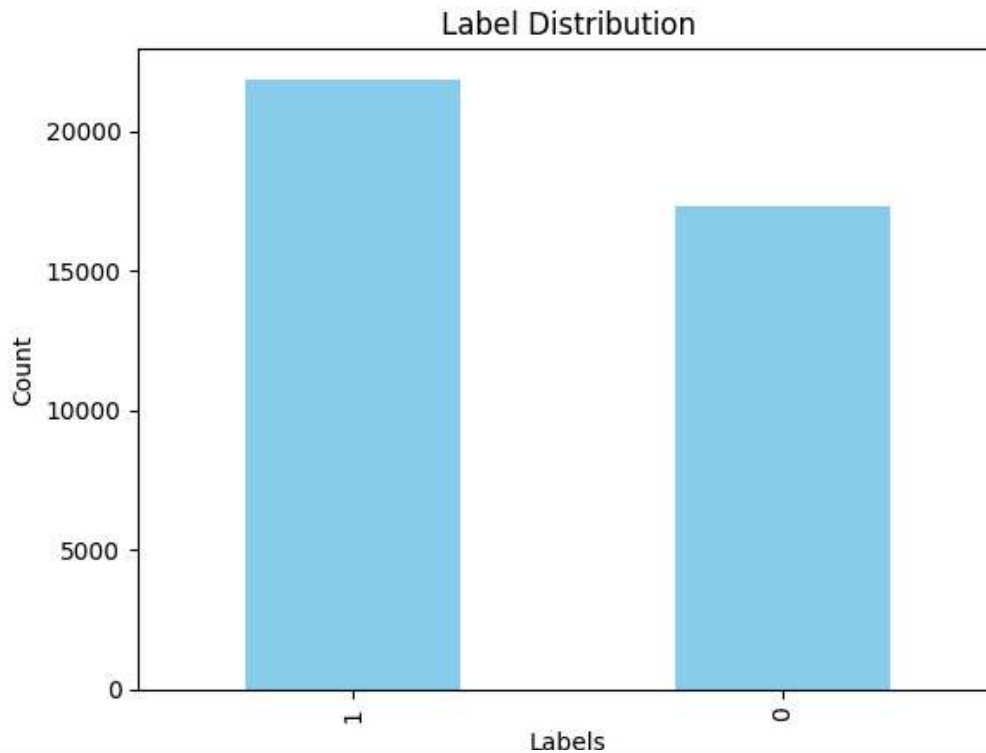
```
274078
```

```
1 import matplotlib.pyplot as plt
2
3 df['label'].value_counts().plot(kind='bar', color='skyblue')
4 plt.title('Label Distribution')
5 plt.xlabel('Labels')
6 plt.ylabel('Count')
7 plt.show()
8
```

⇥

## Label Distribution



```
1 print(df['label'].value_counts())
2
```

⇥  label
     1    21842
     0    17312
   Name: count, dtype: int64

```
 1 import torch
 2 import torch.nn as nn
 3 from torch.utils.data import Dataset, DataLoader
 4 from transformers import BertTokenizer, BertModel, AdamW
 5 import pandas as pd
 6 from sklearn.model_selection import train_test_split
 7 from sklearn.metrics import accuracy_score, f1_score
 8 import re
 9 import numpy as np
10 import gc
11
12 # ----- Data Preprocessing Functions -----
13 def clean_email_text(text):
14     """Clean email text by removing unwanted patterns and normalizing"""
15     if pd.isna(text):
16         return ""
17     text = str(text)
18     text = re.sub(r'From:.*\n|To:.*\n|Subject:.*\n|Date:.*\n', '', text)
19     text = re.sub(r'http[s]?://(?:[a-zA-Z]|[0-9]|[$-_@.&+]|[!*\\(\\),]|(?:%[0-9a-fA-F][0-9a-fA-F]))+',
20     text = re.sub(r'[^a-zA-Z0-9\s.,!?]', '', text)
21     text = re.sub(r'\s+', ' ', text).strip()
22     return text.lower()
23
24 def preprocess_dataframe(df):
25     """Preprocess the entire dataframe"""
```

```python
26      df['combined_text'] = df['subject'].fillna('') + ' ' + df['body'].fillna('')
27      df['combined_text'] = df['combined_text'].apply(clean_email_text)
28      df['label'] = pd.to_numeric(df['label'], errors='coerce').fillna(0).astype(int)
29      return df
30
31  # ----- 1. Prepare the Dataset -----
32  class PhishingEmailDataset(Dataset):
33      def __init__(self, texts, labels, tokenizer, max_length=256):
34          self.texts = texts
35          self.labels = labels
36          self.tokenizer = tokenizer
37          self.max_length = max_length
38
39      def __len__(self):
40          return len(self.texts)
41
42      def __getitem__(self, idx):
43          text = str(self.texts[idx])
44          label = self.labels[idx]
45          if len(text.strip()) == 0:
46              text = "[EMPTY_EMAIL]"
47          encoding = self.tokenizer.encode_plus(
48              text,
49              add_special_tokens=True,
50              max_length=self.max_length,
51              truncation=True,
52              padding='max_length',
53              return_token_type_ids=False,
54              return_attention_mask=True,
55              return_tensors='pt'
56          )
57          return {
58              'input_ids': encoding['input_ids'].flatten(),
59              'attention_mask': encoding['attention_mask'].flatten(),
60              'label': torch.tensor(label, dtype=torch.long)
61          }
62
63  # ----- 2. Define the Hybrid BERT-GRU Model -----
64  class BertGRUClassifier(nn.Module):
65      def __init__(self, n_classes, dropout_rate=0.3, gru_hidden_size=128, num_gru_layers=1):
66          super(BertGRUClassifier, self).__init__()
67          self.bert = BertModel.from_pretrained('bert-base-uncased')
68          self.gru = nn.GRU(input_size=self.bert.config.hidden_size,
69                            hidden_size=gru_hidden_size,
70                            num_layers=num_gru_layers,
71                            batch_first=True,
72                            bidirectional=True)
73          self.dropout = nn.Dropout(dropout_rate)
74          self.out = nn.Linear(gru_hidden_size * 2, n_classes)
75
76      def forward(self, input_ids, attention_mask):
77          bert_outputs = self.bert(input_ids=input_ids, attention_mask=attention_mask)
78          sequence_output = bert_outputs.last_hidden_state
79          gru_output, _ = self.gru(sequence_output)
80          pooled_output = torch.mean(gru_output, dim=1)
81          pooled_output = self.dropout(pooled_output)
82          logits = self.out(pooled_output)
83          return logits
84
```

```python
85 # ----- 3. Load and Preprocess Dataset -----
86 df = pd.read_csv("/content/CEAS_08[1].csv", escapechar='\\')
87 print("Original data shape:", df.shape)
88
89 df_processed = preprocess_dataframe(df)
90 print("Processed data shape:", df_processed.shape)
91 print("Label distribution:", df_processed['label'].value_counts())
92
93 train_texts, val_texts, train_labels, val_labels = train_test_split(
94     df_processed['combined_text'],
95     df_processed['label'],
96     test_size=0.2,
97     random_state=42,
98     stratify=df_processed['label']
99 )
100
101 tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
102 train_dataset = PhishingEmailDataset(train_texts.tolist(), train_labels.tolist(), tokenizer)
103 val_dataset = PhishingEmailDataset(val_texts.tolist(), val_labels.tolist(), tokenizer)
104
105 batch_size = 8  # Reduced batch size to prevent memory overload
106 train_loader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True, num_workers=2)
107 val_loader = DataLoader(val_dataset, batch_size=batch_size, num_workers=2)
108
109 # ----- 4. Training Setup -----
110 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
111 model = BertGRUClassifier(n_classes=2).to(device)
112 optimizer = AdamW(model.parameters(), lr=2e-5)
113 criterion = nn.CrossEntropyLoss()
114
115 def train_epoch(model, data_loader, optimizer, criterion, device):
116     model.train()
117     total_loss = 0
118     for i, batch in enumerate(data_loader):
119         input_ids = batch['input_ids'].to(device)
120         attention_mask = batch['attention_mask'].to(device)
121         labels = batch['label'].to(device)
122
123         optimizer.zero_grad()
124         outputs = model(input_ids=input_ids, attention_mask=attention_mask)
125         loss = criterion(outputs, labels)
126
127         # Gradient clipping to prevent exploding gradients
128         torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
129
130         loss.backward()
131         optimizer.step()
132
133         total_loss += loss.item()
134
135         # Clear memory
136         del input_ids, attention_mask, labels, outputs, loss
137         torch.cuda.empty_cache()
138
139         if i % 10 == 0:  # Print progress every 10 batches
140             print(f"Batch {i}/{len(data_loader)} processed")
141
142     return total_loss / len(data_loader)
143
```

```python
144 def eval_model(model, data_loader, device):
145     model.eval()
146     predictions = []
147     true_labels = []
148     with torch.no_grad():
149         for batch in data_loader:
150             input_ids = batch['input_ids'].to(device)
151             attention_mask = batch['attention_mask'].to(device)
152             labels = batch['label'].to(device)
153             outputs = model(input_ids=input_ids, attention_mask=attention_mask)
154             _, preds = torch.max(outputs, dim=1)
155             predictions.extend(preds.cpu().numpy())
156             true_labels.extend(labels.cpu().numpy())
157
158             # Clear memory
159             del input_ids, attention_mask, labels, outputs
160             torch.cuda.empty_cache()
161
162     return accuracy_score(true_labels, predictions), f1_score(true_labels, predictions, average='weigh
163
164 # ----- 5. Training Loop -----
165 epochs = 10
166 for epoch in range(epochs):
167     print(f"Starting Epoch {epoch+1}/{epochs}")
168     train_loss = train_epoch(model, train_loader, optimizer, criterion, device)
169     val_accuracy, val_f1 = eval_model(model, val_loader, device)
170     print(f"Epoch {epoch+1}/{epochs} | Train Loss: {train_loss:.4f} | Val Acc: {val_accuracy:.4f} | Va
171
172     # Garbage collection
173     gc.collect()
174     torch.cuda.empty_cache()
175
176 print("Training completed!")
```

```
Processed data shape: (39154, 8)
Label distribution: label
1    21842
0    17312
Name: count, dtype: int64
```

tokenizer_config.json: 100%                                  48.0/48.0 [00:00<00:00, 3.81kB/s]

vocab.txt: 100%                                              232k/232k [00:00<00:00, 2.86MB/s]

tokenizer.json: 100%                                         466k/466k [00:00<00:00, 4.41MB/s]

config.json: 100%                                            570/570 [00:00<00:00, 45.0kB/s]

model.safetensors: 100%                                      440M/440M [00:02<00:00, 228MB/s]

```
Starting Epoch 1/10
/usr/local/lib/python3.11/dist-packages/transformers/optimization.py:640: FutureWarning: This implem
  warnings.warn(
Batch 0/3916 processed
Batch 10/3916 processed
Batch 20/3916 processed
Batch 30/3916 processed
Batch 40/3916 processed
Batch 50/3916 processed
Batch 60/3916 processed
Batch 70/3916 processed
Batch 80/3916 processed
Batch 90/3916 processed
Batch 100/3916 processed
Batch 110/3916 processed
Batch 120/3916 processed
Batch 130/3916 processed
Batch 140/3916 processed
Batch 150/3916 processed
Batch 160/3916 processed
Batch 170/3916 processed
Batch 180/3916 processed
Batch 190/3916 processed
Batch 200/3916 processed
Batch 210/3916 processed
Batch 220/3916 processed
Batch 230/3916 processed
Batch 240/3916 processed
Batch 250/3916 processed
Batch 260/3916 processed
Batch 270/3916 processed
Batch 280/3916 processed
Batch 290/3916 processed
Batch 300/3916 processed
Batch 310/3916 processed
Batch 320/3916 processed
Batch 330/3916 processed
Batch 340/3916 processed
Batch 350/3916 processed
Batch 360/3916 processed
Batch 370/3916 processed
Batch 380/3916 processed
Batch 390/3916 processed
Batch 400/3916 processed
Batch 410/3916 processed
Batch 420/3916 processed
Batch 430/3916 processed
Batch 440/3916 processed
Batch 450/3916 processed
Batch 460/3916 processed
```

```
Batch 470/3916 processed
Batch 480/3916 processed
Batch 490/3916 processed
Batch 500/3916 processed
Batch 510/3916 processed
Batch 520/3916 processed
Batch 530/3916 processed
Batch 540/3916 processed
Batch 550/3916 processed
Batch 560/3916 processed
Batch 570/3916 processed
Batch 580/3916 processed
Batch 590/3916 processed
Batch 600/3916 processed
Batch 610/3916 processed
Batch 620/3916 processed
Batch 630/3916 processed
Batch 640/3916 processed
Batch 650/3916 processed
Batch 660/3916 processed
Batch 670/3916 processed
Batch 680/3916 processed
Batch 690/3916 processed
Batch 700/3916 processed
Batch 710/3916 processed
Batch 720/3916 processed
Batch 730/3916 processed
Batch 740/3916 processed
Batch 750/3916 processed
Batch 760/3916 processed
Batch 770/3916 processed
Batch 780/3916 processed
Batch 790/3916 processed
Batch 800/3916 processed
Batch 810/3916 processed
Batch 820/3916 processed
Batch 830/3916 processed
Batch 840/3916 processed
Batch 850/3916 processed
Batch 860/3916 processed
Batch 870/3916 processed
Batch 880/3916 processed
Batch 890/3916 processed
Batch 900/3916 processed
Batch 910/3916 processed
Batch 920/3916 processed
Batch 930/3916 processed
Batch 940/3916 processed
Batch 950/3916 processed
Batch 960/3916 processed
Batch 970/3916 processed
Batch 980/3916 processed
Batch 990/3916 processed
Batch 1000/3916 processed
Batch 1010/3916 processed
Batch 1020/3916 processed
Batch 1030/3916 processed
Batch 1040/3916 processed
Batch 1050/3916 processed
Batch 1060/3916 processed
Batch 1070/3916 processed
Batch 1080/3916 processed
Batch 1090/3916 processed
Batch 1100/3916 processed
Batch 1110/3916 processed
```

```
Batch 1120/3916 processed
Batch 1130/3916 processed
Batch 1140/3916 processed
Batch 1150/3916 processed
Batch 1160/3916 processed
Batch 1170/3916 processed
Batch 1180/3916 processed
Batch 1190/3916 processed
Batch 1200/3916 processed
Batch 1210/3916 processed
Batch 1220/3916 processed
Batch 1230/3916 processed
Batch 1240/3916 processed
Batch 1250/3916 processed
Batch 1260/3916 processed
Batch 1270/3916 processed
Batch 1280/3916 processed
Batch 1290/3916 processed
Batch 1300/3916 processed
Batch 1310/3916 processed
Batch 1320/3916 processed
Batch 1330/3916 processed
Batch 1340/3916 processed
Batch 1350/3916 processed
Batch 1360/3916 processed
Batch 1370/3916 processed
Batch 1380/3916 processed
Batch 1390/3916 processed
Batch 1400/3916 processed
Batch 1410/3916 processed
Batch 1420/3916 processed
Batch 1430/3916 processed
Batch 1440/3916 processed
Batch 1450/3916 processed
Batch 1460/3916 processed
Batch 1470/3916 processed
Batch 1480/3916 processed
Batch 1490/3916 processed
Batch 1500/3916 processed
Batch 1510/3916 processed
Batch 1520/3916 processed
Batch 1530/3916 processed
Batch 1540/3916 processed
Batch 1550/3916 processed
Batch 1560/3916 processed
Batch 1570/3916 processed
Batch 1580/3916 processed
Batch 1590/3916 processed
Batch 1600/3916 processed
Batch 1610/3916 processed
Batch 1620/3916 processed
Batch 1630/3916 processed
Batch 1640/3916 processed
Batch 1650/3916 processed
Batch 1660/3916 processed
Batch 1670/3916 processed
Batch 1680/3916 processed
Batch 1690/3916 processed
Batch 1700/3916 processed
Batch 1710/3916 processed
Batch 1720/3916 processed
Batch 1730/3916 processed
Batch 1740/3916 processed
Batch 1750/3916 processed
Batch 1760/3916 processed
Batch 1770/3916 processed
```

```
Batch 1770/3916 processed
Batch 1780/3916 processed
Batch 1790/3916 processed
Batch 1800/3916 processed
Batch 1810/3916 processed
Batch 1820/3916 processed
Batch 1830/3916 processed
Batch 1840/3916 processed
Batch 1850/3916 processed
Batch 1860/3916 processed
Batch 1870/3916 processed
Batch 1880/3916 processed
Batch 1890/3916 processed
Batch 1900/3916 processed
Batch 1910/3916 processed
Batch 1920/3916 processed
Batch 1930/3916 processed
Batch 1940/3916 processed
Batch 1950/3916 processed
Batch 1960/3916 processed
Batch 1970/3916 processed
Batch 1980/3916 processed
Batch 1990/3916 processed
Batch 2000/3916 processed
Batch 2010/3916 processed
Batch 2020/3916 processed
Batch 2030/3916 processed
Batch 2040/3916 processed
Batch 2050/3916 processed
Batch 2060/3916 processed
Batch 2070/3916 processed
Batch 2080/3916 processed
Batch 2090/3916 processed
Batch 2100/3916 processed
Batch 2110/3916 processed
Batch 2120/3916 processed
Batch 2130/3916 processed
Batch 2140/3916 processed
Batch 2150/3916 processed
Batch 2160/3916 processed
Batch 2170/3916 processed
Batch 2180/3916 processed
Batch 2190/3916 processed
Batch 2200/3916 processed
Batch 2210/3916 processed
Batch 2220/3916 processed
Batch 2230/3916 processed
Batch 2240/3916 processed
Batch 2250/3916 processed
Batch 2260/3916 processed
Batch 2270/3916 processed
Batch 2280/3916 processed
Batch 2290/3916 processed
Batch 2300/3916 processed
Batch 2310/3916 processed
Batch 2320/3916 processed
Batch 2330/3916 processed
Batch 2340/3916 processed
Batch 2350/3916 processed
Batch 2360/3916 processed
Batch 2370/3916 processed
Batch 2380/3916 processed
Batch 2390/3916 processed
Batch 2400/3916 processed
Batch 2410/3916 processed
Batch 2420/3916 processed
```

```
Batch 2430/3916 processed
Batch 2440/3916 processed
Batch 2450/3916 processed
Batch 2460/3916 processed
Batch 2470/3916 processed
Batch 2480/3916 processed
Batch 2490/3916 processed
Batch 2500/3916 processed
Batch 2510/3916 processed
Batch 2520/3916 processed
Batch 2530/3916 processed
Batch 2540/3916 processed
Batch 2550/3916 processed
Batch 2560/3916 processed
Batch 2570/3916 processed
Batch 2580/3916 processed
Batch 2590/3916 processed
Batch 2600/3916 processed
Batch 2610/3916 processed
Batch 2620/3916 processed
Batch 2630/3916 processed
Batch 2640/3916 processed
Batch 2650/3916 processed
Batch 2660/3916 processed
Batch 2670/3916 processed
Batch 2680/3916 processed
Batch 2690/3916 processed
Batch 2700/3916 processed
Batch 2710/3916 processed
Batch 2720/3916 processed
Batch 2730/3916 processed
Batch 2740/3916 processed
Batch 2750/3916 processed
Batch 2760/3916 processed
Batch 2770/3916 processed
Batch 2780/3916 processed
Batch 2790/3916 processed
Batch 2800/3916 processed
Batch 2810/3916 processed
Batch 2820/3916 processed
Batch 2830/3916 processed
Batch 2840/3916 processed
Batch 2850/3916 processed
Batch 2860/3916 processed
Batch 2870/3916 processed
Batch 2880/3916 processed
Batch 2890/3916 processed
Batch 2900/3916 processed
Batch 2910/3916 processed
Batch 2920/3916 processed
Batch 2930/3916 processed
Batch 2940/3916 processed
Batch 2950/3916 processed
Batch 2960/3916 processed
Batch 2970/3916 processed
Batch 2980/3916 processed
Batch 2990/3916 processed
Batch 3000/3916 processed
Batch 3010/3916 processed
Batch 3020/3916 processed
Batch 3030/3916 processed
Batch 3040/3916 processed
Batch 3050/3916 processed
Batch 3060/3916 processed
Batch 3070/3916 processed
Batch 3080/3916 processed
```

```
Batch 3080/3916 processed
Batch 3090/3916 processed
Batch 3100/3916 processed
Batch 3110/3916 processed
Batch 3120/3916 processed
Batch 3130/3916 processed
Batch 3140/3916 processed
Batch 3150/3916 processed
Batch 3160/3916 processed
Batch 3170/3916 processed
Batch 3180/3916 processed
Batch 3190/3916 processed
Batch 3200/3916 processed
Batch 3210/3916 processed
Batch 3220/3916 processed
Batch 3230/3916 processed
Batch 3240/3916 processed
Batch 3250/3916 processed
Batch 3260/3916 processed
Batch 3270/3916 processed
Batch 3280/3916 processed
Batch 3290/3916 processed
Batch 3300/3916 processed
Batch 3310/3916 processed
Batch 3320/3916 processed
Batch 3330/3916 processed
Batch 3340/3916 processed
Batch 3350/3916 processed
Batch 3360/3916 processed
Batch 3370/3916 processed
Batch 3380/3916 processed
Batch 3390/3916 processed
Batch 3400/3916 processed
Batch 3410/3916 processed
Batch 3420/3916 processed
Batch 3430/3916 processed
Batch 3440/3916 processed
Batch 3450/3916 processed
Batch 3460/3916 processed
Batch 3470/3916 processed
Batch 3480/3916 processed
Batch 3490/3916 processed
Batch 3500/3916 processed
Batch 3510/3916 processed
Batch 3520/3916 processed
Batch 3530/3916 processed
Batch 3540/3916 processed
Batch 3550/3916 processed
Batch 3560/3916 processed
Batch 3570/3916 processed
Batch 3580/3916 processed
Batch 3590/3916 processed
Batch 3600/3916 processed
Batch 3610/3916 processed
Batch 3620/3916 processed
Batch 3630/3916 processed
Batch 3640/3916 processed
Batch 3650/3916 processed
Batch 3660/3916 processed
Batch 3670/3916 processed
Batch 3680/3916 processed
Batch 3690/3916 processed
Batch 3700/3916 processed
Batch 3710/3916 processed
Batch 3720/3916 processed
Batch 3730/3916 processed
```

```
Batch 3740/3916 processed
Batch 3750/3916 processed
Batch 3760/3916 processed
Batch 3770/3916 processed
Batch 3780/3916 processed
Batch 3790/3916 processed
Batch 3800/3916 processed
Batch 3810/3916 processed
Batch 3820/3916 processed
Batch 3830/3916 processed
Batch 3840/3916 processed
Batch 3850/3916 processed
Batch 3860/3916 processed
Batch 3870/3916 processed
Batch 3880/3916 processed
Batch 3890/3916 processed
Batch 3900/3916 processed
Batch 3910/3916 processed
Epoch 1/10 | Train Loss: 0.0288 | Val Acc: 0.9980 | Val F1: 0.9980
Starting Epoch 2/10
Batch 0/3916 processed
Batch 10/3916 processed
Batch 20/3916 processed
Batch 30/3916 processed
Batch 40/3916 processed
Batch 50/3916 processed
Batch 60/3916 processed
Batch 70/3916 processed
Batch 80/3916 processed
Batch 90/3916 processed
Batch 100/3916 processed
Batch 110/3916 processed
Batch 120/3916 processed
Batch 130/3916 processed
Batch 140/3916 processed
Batch 150/3916 processed
Batch 160/3916 processed
Batch 170/3916 processed
Batch 180/3916 processed
Batch 190/3916 processed
Batch 200/3916 processed
Batch 210/3916 processed
Batch 220/3916 processed
Batch 230/3916 processed
Batch 240/3916 processed
Batch 250/3916 processed
Batch 260/3916 processed
Batch 270/3916 processed
Batch 280/3916 processed
Batch 290/3916 processed
Batch 300/3916 processed
Batch 310/3916 processed
Batch 320/3916 processed
Batch 330/3916 processed
Batch 340/3916 processed
Batch 350/3916 processed
Batch 360/3916 processed
Batch 370/3916 processed
Batch 380/3916 processed
Batch 390/3916 processed
Batch 400/3916 processed
Batch 410/3916 processed
Batch 420/3916 processed
Batch 430/3916 processed
Batch 440/3916 processed
Batch 450/3916 processed
```