```python
1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import math
5 import numpy as np
6 import torch.cuda.amp as amp
7 from tqdm import tqdm
8 import matplotlib.pyplot as plt
9 from torch.utils.data import DataLoader, TensorDataset
```

```python
 1 # Training Configuration
 2 epochs = 5
 3 batch_size = 1024
 4 lr = 3e-4
 5 weight_decay = 0.01
 6 device = "cuda"
 7 checkpoint_filepath = None  # Set to a path if you want to load a checkpoint
 8 save_dir = "checkpoints"
 9 import os
10 os.makedirs(save_dir, exist_ok=True)
```

## ˅ Model code

```python
 1 class PatchEmbedding(nn.Module):
 2     """
 3     Module that converts image patches to embeddings for Vision Transformer.
 4     """
 5     def __init__(self,
 6                  image_size: tuple = (64, 72),
 7                  patch_size: int = 8,
 8                  in_channels: int = 3,
 9                  embedding_dim: int = 1024):
10         super().__init__()
11         self.image_size = image_size
12         self.patch_size = patch_size
13         self.in_channels = in_channels
14
15         # Calculate number of patches
16         self.num_patches = (image_size[0] // patch_size) * (image_size[1] // p
17
```

```
18          # Create projection for converting patches to embeddings
19          self.projection = nn.Conv2d(
20              in_channels=in_channels,
21              out_channels=embedding_dim,
22              kernel_size=patch_size,
23              stride=patch_size
24          )
25
26          # CLS token embedding
27          self.cls_token = nn.Parameter(torch.zeros(1, 1, embedding_dim))
28
29          # Positional embedding (Normal distribution initialization of value)
30          self.positions = nn.Parameter(torch.zeros(1, self.num_patches + 1, emb
31          nn.init.trunc_normal_(self.positions, std=0.02)
32
33      def forward(self, x: torch.Tensor) -> torch.Tensor:
34          batch_size = x.shape[0]
35
36          # Convert image to patches and project to embedding dimension
37          # x shape: [batch_size, channels, height, width]
38          x = self.projection(x)
39          # x shape: [batch_size, embedding_dim, height/patch_size, width/patch_
40
41          # Flatten patches to sequence
42          x = x.flatten(2).transpose(1, 2)
43          # x shape: [batch_size, num_patches, embedding_dim]
44
45          # Add CLS token
46          cls_tokens = self.cls_token.expand(batch_size, -1, -1)
47          x = torch.cat((cls_tokens, x), dim=1)
48
49          # Add positional embeddings
50          x = x + self.positions
51
52          return x
53
54
55 class VisionAttention(nn.Module):
56     def __init__(self,
57                  hidden_dim: int,
58                  head_dim: int,
59                  q_head: int,
60                  kv_head: int,
61                  lora_rank: int = 16):
62         super().__init__()
```

```
63              self.head_dim = head_dim
64              self.q_head = q_head
65              self.kv_head = kv_head
66              self.qkv = nn.Linear(hidden_dim, (q_head+kv_head*2)*head_dim)
67              self.o = nn.Linear(q_head*head_dim, hidden_dim)
68              self.scaler = 1/math.sqrt(head_dim)
69              self.lora_qkv_a = nn.Linear(hidden_dim, lora_rank)
70              self.lora_qkv_b = nn.Linear(lora_rank, (q_head+kv_head*2)*head_dim)
71              self.lora_o_a = nn.Linear(q_head*head_dim, lora_rank)
72              self.lora_o_b = nn.Linear(lora_rank, hidden_dim)
73
74              if q_head != kv_head:
75                  # If we are using multi query attention
76                  assert q_head % kv_head == 0
77                  self.multi_query_attention = True
78                  self.q_kv_scale = q_head//kv_head
79              else:
80                  self.multi_query_attention = False
81
82          def forward(self, tensor: torch.Tensor, attention_mask: torch.Tensor = Nor
83              batch_size, seq_len, hid_dim = tensor.shape
84
85              qkv_tensor = self.qkv(tensor)
86              if fine_tuning:
87                  lora_tensor = self.lora_qkv_a(tensor)
88                  lora_tensor = self.lora_qkv_b(lora_tensor)
89                  qkv_tensor = lora_tensor + qkv_tensor
90              query, key, value = qkv_tensor.split([self.head_dim*self.q_head, self.
91
92              query = query.view(batch_size, seq_len, self.q_head, self.head_dim)
93              key = key.view(batch_size, seq_len, self.kv_head, self.head_dim)
94              value = value.view(batch_size, seq_len, self.kv_head, self.head_dim)
95
96              if self.multi_query_attention:
97                  # If we are using multi query attention, duplicate key value heads
98                  key = torch.repeat_interleave(key, self.q_kv_scale, dim=-2)
99                  value = torch.repeat_interleave(value, self.q_kv_scale, dim=-2)
100
101             # Switch to batch_size, head, seq_len, head_dim
102             query = query.transpose(1, 2)
103             key = key.transpose(1, 2)
104             value = value.transpose(1, 2)
105
106             # Classic self attention
107             attention_raw = torch.matmul(query, key.transpose(2, 3))
```

```python
108            attention_scaled = attention_raw * self.scaler
109            if attention_mask != None:
110                attention_scaled += attention_mask
111            attention_score = torch.softmax(attention_scaled, dim=-1)
112            value = torch.matmul(attention_score, value)
113
114            # Reshape back to batch_size, seq_len, hid_dim
115            value = value.transpose(1, 2).contiguous()
116            value = value.view(batch_size, seq_len, hid_dim)
117
118            # Output layer
119            output = self.o(value)
120            if fine_tuning:
121                lora_tensor = self.lora_o_a(value)
122                lora_tensor = self.lora_o_b(lora_tensor)
123                output = lora_tensor + output
124
125            return output
126
127
128 class FeedForward(nn.Module):
129     def __init__(self,
130                  hidden_size: int,
131                  expansion_factor: int = 4,
132                  dropout_ratio: float = 0.1,
133                  lora_rank: int = 16):
134         super().__init__()
135         self.gate_and_up = nn.Linear(hidden_size, hidden_size * expansion_fact
136         self.down = nn.Linear(hidden_size * expansion_factor, hidden_size)
137         self.dropout = nn.Dropout(p=dropout_ratio)
138         self.lora_gate_and_up_a = nn.Linear(hidden_size, lora_rank)
139         self.lora_gate_and_up_b = nn.Linear(lora_rank, hidden_size * expansion
140         self.lora_down_a = nn.Linear(hidden_size * expansion_factor, lora_ranl
141         self.lora_down_b = nn.Linear(lora_rank, hidden_size)
142
143     def forward(self, tensor: torch.Tensor, fine_tuning: bool = False) -> tor
144         gate_and_up = self.gate_and_up(tensor)
145         if fine_tuning:
146             lora_tensor = self.lora_gate_and_up_a(tensor)
147             lora_tensor = self.lora_gate_and_up_b(lora_tensor)
148             gate_and_up = gate_and_up + lora_tensor
149         gate, up = gate_and_up.chunk(chunks=2, dim=-1)
150         gate = F.gelu(gate, approximate="tanh")
151         tensor = gate * up
152         tensor = self.dropout(tensor)
153         down_tensor = self.down(tensor)
```

```python
153         down_tensor = self.down(tensor)
154         if fine_tuning:
155             lora_tensor = self.lora_down_a(tensor)
156             lora_tensor = self.lora_down_b(lora_tensor)
157             down_tensor = down_tensor + lora_tensor
158         return down_tensor
159
160
161 class MOE(nn.Module):
162     def __init__(self, hidden_size: int, device: str, num_experts: int = 8, e
163         super().__init__()
164         self.gate = nn.Linear(hidden_size, num_experts)
165         self.num_experts = num_experts
166         self.device = device
167         self.experts = nn.ModuleList([FeedForward(hidden_size, expansion_facto
168
169     def forward(self, tensor: torch.Tensor, fine_tuning: bool = False) -> tup
170         # Flatten for better manipulation, this is ok because tokens are indep
171         batch_size, seq_len, hidden_size = tensor.shape
172         flat_tensor = tensor.reshape(batch_size * seq_len, hidden_size)
173
174         # Pass through the gating network and select experts
175         tensor = self.gate(flat_tensor)
176         tensor = F.softmax(tensor, dim=-1)
177
178         # The output of this step is a tensor of shape [batch_size * seq_len,
179         value_tensor, index_tensor = tensor.topk(k=2, dim=-1)
180
181         # Find the load balancing loss
182         counts = torch.bincount(index_tensor[:, 0], minlength=self.num_experts
183         frequencies = counts.float() / (batch_size * seq_len) # This is the ha
184         probability = tensor.mean(0) # This is the soft probability
185         load_balancing_loss = (probability * frequencies).mean() * float(self.
186
187         # Normalize top1 and top2 score
188         top_expert_score = value_tensor[:, 0]
189         second_expert_score = value_tensor[:, 1]
190         total_score = top_expert_score + second_expert_score
191         top_expert_score = top_expert_score / total_score
192         second_expert_score = second_expert_score / total_score
193
194         # Split into top 2 experts
195         split_tensors = torch.split(index_tensor, 1, dim=-1)
196         top_expert, second_expert = split_tensors[0], split_tensors[1]
197         indices = torch.arange(batch_size * seq_len).unsqueeze(-1).to(self.dev
198         top_expert = torch.cat((indices, top_expert), dim=-1)
```

```python
198         top_expert = torch.cat((indices, top_expert), dim=-1)
199         second_expert = torch.cat((indices, second_expert), dim=-1)
200
201         # Sort based on expert selection
202         top_expert = top_expert[top_expert[:,1].argsort()]
203         second_expert = second_expert[second_expert[:,1].argsort()]
204
205         # Count how many tokens goes to each expert
206         top_expert_counts = torch.zeros(self.num_experts, dtype=int)
207         for i in range(self.num_experts):
208             top_expert_counts[i] = (top_expert[:,1] == i).sum()
209         top_expert_counts = top_expert_counts.tolist()
210
211         second_expert_counts = torch.zeros(self.num_experts, dtype=int)
212         for i in range(self.num_experts):
213             second_expert_counts[i] = (second_expert[:,1] == i).sum()
214         second_expert_counts = second_expert_counts.tolist()
215
216         # Split input tokens for each expert
217         top_expert_tokens = flat_tensor[top_expert[:,0]]
218         second_expert_tokens = flat_tensor[second_expert[:,0]]
219
220         # Split into a list of tensors, element i tensor is for ith expert.
221         top_expert_tokens = torch.split(top_expert_tokens, top_expert_counts,
222         second_expert_tokens = torch.split(second_expert_tokens, second_exper
223
224         # Input into each expert and obtain results in a list
225         top_expert_outputs = [self.experts[i](top_expert_tokens[i], fine_tunin
226         second_expert_outputs = [self.experts[i](second_expert_tokens[i], fine
227
228         # Combine outputs
229         top_expert_outputs = torch.cat(top_expert_outputs, dim=0)
230         second_expert_outputs = torch.cat(second_expert_outputs, dim=0)
231
232         # Re-index the output back to original token order
233         flat_top_expert_tensor = torch.zeros_like(flat_tensor, dtype=torch.flo
234         flat_top_expert_tensor.index_copy_(0, top_expert[:, 0], top_expert_out
235
236         flat_second_expert_tensor = torch.zeros_like(flat_tensor, dtype=torch.
237         flat_second_expert_tensor.index_copy_(0, second_expert[:, 0], second_e
238
239         # Find final output tensor based on weight between top and second expe
240         final_tensor = top_expert_score.unsqueeze(-1) * flat_top_expert_tensor
241
242         # Reshape to original [batch_size, seq_len, hidden_size]
243         final_tensor = final_tensor.reshape(batch_size, seq_len, hidden_size)
```

```python
243         final_tensor = final_tensor.reshape(batch_size, seq_len, hidden_size)
244
245         return final_tensor, load_balancing_loss
246
247
248 class VisionLayer(nn.Module):
249     def __init__(self,
250                  hidden_dim: int,
251                  head_dim: int,
252                  q_head: int,
253                  kv_head: int,
254                  device: str,
255                  expansion_factor: int = 4,
256                  dropout_ratio: float = 0.1,
257                  use_moe: bool = False,
258                  num_experts: int = 8,
259                  lora_rank: int = 16):
260         super().__init__()
261         self.use_moe = use_moe
262         self.device = device
263
264         self.norm1 = nn.LayerNorm(hidden_dim)
265         self.attention = VisionAttention(hidden_dim, head_dim, q_head, kv_head
266
267         self.norm2 = nn.LayerNorm(hidden_dim)
268         if self.use_moe:
269             self.moe = MOE(hidden_dim, device, num_experts=num_experts, expans
270                            dropout_ratio=dropout_ratio, lora_rank=lora_rank)
271         else:
272             self.ffn = FeedForward(hidden_dim, expansion_factor=expansion_fact
273                                    lora_rank=lora_rank)
274
275     def forward(self, tensor: torch.Tensor, attention_mask: torch.Tensor = No
276         skip_connection = tensor
277         tensor = self.norm1(tensor)
278         tensor = self.attention(tensor, attention_mask=attention_mask, fine_tu
279         tensor += skip_connection
280
281         skip_connection = tensor
282         tensor = self.norm2(tensor)
283         if self.use_moe:
284             tensor, load_balancing_loss = self.moe(tensor, fine_tuning=fine_tu
285         else:
286             tensor = self.ffn(tensor, fine_tuning=fine_tuning)
287             load_balancing_loss = torch.tensor(0.0, dtype=tensor.dtype, device
288
```

```python
289              tensor += skip_connection
290
291              return tensor, load_balancing_loss
292
293
294  class VisionTransformer(nn.Module):
295      def __init__(self,
296                      image_size: tuple,
297                      num_classes: int = 1,
298                      patch_size: int = 8,
299                      in_channels: int = 3,
300                      num_layer: int = 3,
301                      hidden_dim: int = 1024,
302                      expansion_factor: int = 8,
303                      head_dim: int = 64,
304                      q_head: int = 16,
305                      kv_head: int = 4,
306                      dropout_ratio: float = 0.1,
307                      use_moe: bool = True,
308                      num_experts: int = 8,
309                      load_balancing_loss_weight: float = 1e-2,
310                      fine_tuning: bool = False,
311                      lora_rank: int = 16):
312          super().__init__()
313          self.device = device
314          self.num_layer = num_layer
315          self.load_balancing_loss_weight = load_balancing_loss_weight
316          self.fine_tuning = fine_tuning
317
318          # Patch embedding
319          self.patch_embedding = PatchEmbedding(
320              image_size=image_size,
321              patch_size=patch_size,
322              in_channels=in_channels,
323              embedding_dim=hidden_dim
324          )
325
326          # Calculate number of patches (sequence length)
327          self.num_patches = (image_size[0] // patch_size) * (image_size[1] // p
328
329          if q_head == None:
330              q_head = (hidden_dim // head_dim)
331
332          if kv_head == None:
333              kv_head = (hidden_dim // head_dim)
```

```python
334
335            if hidden_dim % (head_dim * q_head) != 0 or hidden_dim % (head_dim * I
336                raise ValueError("Error: hidden_dim or projection_dim (if specifi(
337
338            # Create transformer layers
339            self.transformer = nn.ModuleList()
340            for _ in range(self.num_layer):
341                self.transformer.append(VisionLayer(
342                    hidden_dim, head_dim, q_head, kv_head, device,
343                    expansion_factor=expansion_factor,
344                    dropout_ratio=dropout_ratio,
345                    use_moe=use_moe,
346                    num_experts=num_experts,
347                    lora_rank=lora_rank
348                ))
349            self.output_norm = nn.LayerNorm(hidden_dim)
350
351            # Final classifier head
352            self.classifier = nn.Linear(hidden_dim, num_classes)
353
354    def begin_fine_tunning(self) -> None:
355        self.fine_tuning = True
356        for name, param in self.named_parameters():
357            if "lora" not in name:
358                param.requires_grad = False
359            else:
360                param.requires_grad = True
361
362    def exit_fine_tunning(self) -> None:
363        self.fine_tuning = False
364        for name, param in self.named_parameters():
365            if "positions" in name:
366                param.requires_grad = False
367            else:
368                param.requires_grad = True
369
370    def forward(self, x: torch.Tensor) -> tuple[torch.Tensor, torch.Tensor]:
371        # Handle input shape
372        if len(x.shape) == 3:  # [batch_size, 64, 72]
373            # Reshape to [batch_size, channels, height, width]
374            # Assuming the input is grayscale (1 channel)
375            batch_size, height, width = x.shape
376            x = x.unsqueeze(1)  # Add channel dimension [batch_size, 1, 64, 7:
377
378        # Apply patch embedding
```

```
379            x = self.patch_embedding(x)
380
381            # Track load-balancing across layers (only if MoE is used)
382            load_balancing_sum = torch.tensor(0.0, device=self.device)
383
384            # Pass through transformer layers
385            for layer in self.transformer:
386                x, load_balancing_loss = layer(x, fine_tuning=self.fine_tuning)
387                load_balancing_sum += load_balancing_loss
388
389            load_balancing_loss = (load_balancing_sum / self.num_layer) * self.loa
390
391            # Apply output normalization
392            x = self.output_norm(x)
393
394            # Use CLS token for classification
395            x = x[:, 0]  # Take only the CLS token
396
397            # Apply classifier
398            x = self.classifier(x)
399
400            return x, load_balancing_loss
```

## Preprocessing Data

```
 1 data0 = np.load('Run357479_Dataset_iodic.npy')
 2 data1 = np.load('Run355456_Dataset_jqkne.npy')
 3
 4 # Create labels: 0 for class 0 and 1 for class 1
 5 labels_0 = np.zeros((data0.shape[0],), dtype=np.int32)
 6 labels_1 = np.ones((data1.shape[0],), dtype=np.int32)
 7
 8 # Concatenate data and labels
 9 data = np.concatenate([data0, data1], axis=0)  # Shape (20000, 64, 72)
10 labels = np.concatenate([labels_0, labels_1], axis=0)  # Shape (20000,)
11
12 # Shuffle data and labels together
13 indices = np.random.permutation(data.shape[0])
14 data = data[indices]
15 labels = labels[indices]
16
17 # Seperate into training, validation and testing set.
```

```
18 n_total = data.shape[0]
19 n_train = int(n_total * 0.8)
20 n_val = int(n_total * 0.1)
21 n_test = n_total - n_train - n_val
22
23 train_data, val_data, test_data = data[:n_train], data[n_train:n_train + n_va
24 train_labels, val_labels, test_labels = labels[:n_train], labels[n_train:n_tr
25
26 # Compute mean and std from the training data only
27 mean = train_data.mean()
28 std = train_data.std()
29
30 # Apply the same transformation to train, val, test
31 train_data = (train_data - mean) / (std + 1e-7)
32 val_data = (val_data - mean) / (std + 1e-7)
33 test_data = (test_data - mean) / (std + 1e-7)
34
35 # Convert to PyTorch tensors
36 train_data_tensor = torch.tensor(train_data, dtype=torch.float32)
37 val_data_tensor = torch.tensor(val_data, dtype=torch.float32)
38 test_data_tensor = torch.tensor(test_data, dtype=torch.float32)
39
40 train_labels_tensor = torch.tensor(train_labels, dtype=torch.long)
41 val_labels_tensor = torch.tensor(val_labels, dtype=torch.long)
42 test_labels_tensor = torch.tensor(test_labels, dtype=torch.long)
43
44 # Create PyTorch DataLoaders
45 train_dataset = TensorDataset(train_data_tensor, train_labels_tensor.float().
46 val_dataset = TensorDataset(val_data_tensor, val_labels_tensor.float().unsque
47 test_dataset = TensorDataset(test_data_tensor, test_labels_tensor.float().uns
48
49 train_loader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)
50 val_loader = DataLoader(val_dataset, batch_size=batch_size)
51 test_loader = DataLoader(test_dataset, batch_size=batch_size)
```

## ˅ Preparing for training

```python
# Checkpoint loading
def load_checkpoint(model, optimizer, filepath):
    checkpoint = torch.load(filepath)
    model.load_state_dict(checkpoint['model_state_dict'])
    optimizer.load_state_dict(checkpoint['optimizer_state_dict'])
    epoch = checkpoint['epoch']
    validation_loss = checkpoint.get('validation_loss', float('inf'))
    print(f"Loaded checkpoint from epoch {epoch} with validation loss {valida
    return epoch

# Checkpoint saving
def save_checkpoint(model, optimizer, epoch, validation_loss):
    checkpoint = {
        'model_state_dict': model.state_dict(),
        'optimizer_state_dict': optimizer.state_dict(),
        'epoch': epoch,
        'validation_loss': validation_loss
    }
    torch.save(checkpoint, f"{save_dir}/vit_checkpoint_epoch_{epoch}.pt")
    # Save best model separately
    if epoch == 0 or validation_loss < min(loss_valid):
        torch.save(checkpoint, f"{save_dir}/vit_best_model.pt")
        print(f"Saved best model with validation loss: {validation_loss:.6f}"
```

```python
 1  vit = VisionTransformer(
 2      image_size=(64, 72),      # Your input image dimensions
 3      patch_size=8,             # Size of each patch
 4      in_channels=1,            # We technically don't have channel value here.
 5      num_classes=1,            # Number of output classes
 6      num_layer=3,              # Number of transformer layers
 7      hidden_dim=256,          # Hidden dimension
 8      expansion_factor=4,       # Expansion factor for FFN
 9      head_dim=64,              # Dimension of each attention head
10      q_head=4,                 # Number of query heads
11      kv_head=1,                # Number of key/value heads
12      use_moe=True,             # Whether to use Mixture of Experts
13      num_experts=4
14  ).to(device)
15
16  # Load checkpoint if available
17  current_epoch = 0
18  if checkpoint_filepath is not None and checkpoint_filepath != "":
19      current_epoch = load_checkpoint(vit, optimizer, checkpoint_filepath) + 1
20
21  print(f"This model has {sum(p.numel() for p in vit.parameters())} parameters.'
22  print(f"Training on {device}")
```

```
This model has 10774253 parameters.
Training on cuda
```

```python
1  # Binary classification loss since num_classes=1
2  criterion = nn.BCEWithLogitsLoss()
3  optimizer = torch.optim.AdamW(vit.parameters(), lr=lr, weight_decay=weight_de
4  scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, T_max=epoch
5  scaler = amp.GradScaler() # Initialize gradient scaler for mixed precision tr
```

```python
1  # Initialize loss tracking lists
2  loss_train = []
3  loss_valid = []
4  accuracy_train = []
5  accuracy_valid = []
```

## ∨ Training

```python
1  # Training loop
```

```python
1 # Training loop
2 for epoch in range(current_epoch, epochs):
3     print(f"Epoch {epoch+1}/{epochs}")
4
5     # Training phase
6     vit.train()
7     loss_train_epoch = []
8     correct_train = 0
9     total_train = 0
10
11     for inputs, targets in tqdm(train_loader, desc="Training"):
12         inputs = inputs.to(device).float()
13         targets = targets.to(device).float()
14
15         # Forward pass with mixed precision
16         with amp.autocast():
17             outputs, load_balancing_loss = vit(inputs)
18             loss = criterion(outputs, targets) + load_balancing_loss
19
20         # Backward pass with gradient scaling
21         scaler.scale(loss).backward()
22
23         # Gradient clipping
24         scaler.unscale_(optimizer)
25         torch.nn.utils.clip_grad_norm_(vit.parameters(), max_norm=1.0)
26
27         # Optimizer step with scaler
28         scaler.step(optimizer)
29         scaler.update()
30
31         # Zero gradients
32         optimizer.zero_grad()
33
34         # Update scheduler
35         scheduler.step()
36
37         # Record loss
38         loss_train_epoch.append(loss.item())
39
40         # Calculate accuracy
41         predicted = (torch.sigmoid(outputs) > 0.5).float()
42         total_train += targets.size(0)
43         correct_train += (predicted == targets).sum().item()
44
45     # Calculate epoch statistics
```

```
46        epoch_loss = np.mean(loss_train_epoch)
47        epoch_accuracy = 100 * correct_train / total_train
48        loss_train.append(epoch_loss)
49        accuracy_train.append(epoch_accuracy)
50
51        # Validation phase
52        vit.eval()
53        loss_val_epoch = []
54        correct_val = 0
55        total_val = 0
56
57        with torch.no_grad():
58            for inputs, targets in tqdm(val_loader, desc="Validation"):
59                inputs = inputs.to(device).float()
60                targets = targets.to(device).float()
61
62                # Forward pass
63                with amp.autocast():
64                    outputs, load_balancing_loss = vit(inputs)
65                    loss = criterion(outputs, targets) + load_balancing_loss
66
67                # Record loss
68                loss_val_epoch.append(loss.item())
69
70                # Calculate accuracy
71                predicted = (torch.sigmoid(outputs) > 0.5).float()
72                total_val += targets.size(0)
73                correct_val += (predicted == targets).sum().item()
74
75        # Calculate epoch validation statistics
76        epoch_val_loss = np.mean(loss_val_epoch)
77        epoch_val_accuracy = 100 * correct_val / total_val
78        loss_valid.append(epoch_val_loss)
79        accuracy_valid.append(epoch_val_accuracy)
80
81        # Print epoch results
82        print(f"Training — Loss: {epoch_loss:.6f}, Accuracy: {epoch_accuracy:.2f
83        print(f"Validation — Loss: {epoch_val_loss:.6f}, Accuracy: {epoch_val_ac
84
85        # Plot and save training progress
86        plt.figure(figsize=(12, 5))
87
88        plt.subplot(1, 2, 1)
89        plt.plot(loss_train, label="Training loss")
90        plt.plot(loss_valid, label="Validation loss")
```

```python
 91      plt.xlabel("Epoch")
 92      plt.ylabel("Loss")
 93      plt.legend()
 94      plt.title("Training and Validation Loss")
 95
 96      plt.subplot(1, 2, 2)
 97      plt.plot(accuracy_train, label="Training accuracy")
 98      plt.plot(accuracy_valid, label="Validation accuracy")
 99      plt.xlabel("Epoch")
100      plt.ylabel("Accuracy (%)")
101      plt.legend()
102      plt.title("Training and Validation Accuracy")
103
104      plt.tight_layout()
105      plt.savefig(f"{save_dir}/training_progress_epoch_{epoch}.png")
106      plt.close()
107
108  # Save checkpoint
109  save_checkpoint(vit, optimizer, epoch, epoch_val_loss)
```

```
Epoch 1/5
Training: 100%|████████████| 16/16 [00:06<00:00,  2.63it/s]
Validation: 100%|████████████| 2/2 [00:00<00:00,  3.58it/s]
Training — Loss: 0.579171, Accuracy: 74.58%
Validation — Loss: 0.027781, Accuracy: 99.70%
Epoch 2/5
Training: 100%|████████████| 16/16 [00:04<00:00,  3.26it/s]
Validation: 100%|████████████| 2/2 [00:00<00:00,  3.17it/s]
Training — Loss: 0.016199, Accuracy: 99.88%
Validation — Loss: 0.010464, Accuracy: 100.00%
Epoch 3/5
Training: 100%|████████████| 16/16 [00:04<00:00,  3.29it/s]
Validation: 100%|████████████| 2/2 [00:00<00:00,  3.11it/s]
Training — Loss: 0.010416, Accuracy: 100.00%
Validation — Loss: 0.010539, Accuracy: 100.00%
Epoch 4/5
Training: 100%|████████████| 16/16 [00:04<00:00,  3.20it/s]
Validation: 100%|████████████| 2/2 [00:00<00:00,  3.57it/s]
Training — Loss: 0.010518, Accuracy: 100.00%
Validation — Loss: 0.010374, Accuracy: 100.00%
Epoch 5/5
Training: 100%|████████████| 16/16 [00:05<00:00,  3.15it/s]
Validation: 100%|████████████| 2/2 [00:00<00:00,  3.07it/s]
Training — Loss: 0.010164, Accuracy: 100.00%
Validation — Loss: 0.010069, Accuracy: 100.00%
```

```python
  1  import numpy as np
```

```python
2  import matplotlib.pyplot as plt
3  from sklearn.metrics import roc_curve, auc
4
5  # Lists to store true labels and predicted probabilities
6  all_targets = []
7  all_probs = []
8
9  print("\nEvaluating on test set...")
10 vit.eval()
11 test_loss = 0
12 correct = 0
13 total = 0
14
15 with torch.no_grad():
16     for inputs, targets in tqdm(test_loader, desc="Testing"):
17         inputs = inputs.to(device).float()
18         targets = targets.to(device).float()
19
20         with amp.autocast():
21             outputs, load_balancing_loss = vit(inputs)
22             loss = criterion(outputs, targets) + load_balancing_loss
23
24         test_loss += loss.item()
25         # Compute probabilities using sigmoid
26         probabilities = torch.sigmoid(outputs)
27         # Calculate binary predictions for accuracy
28         predicted = (probabilities > 0.5).float()
29         total += targets.size(0)
30         correct += (predicted == targets).sum().item()
31
32         # Append to lists (move to CPU and convert to numpy)
33         all_targets.append(targets.cpu().numpy())
34         all_probs.append(probabilities.cpu().numpy())
35
36 avg_test_loss = test_loss / len(test_loader)
37 test_accuracy = 100 * correct / total
38
39 print(f"Test set – Loss: {avg_test_loss:.6f}, Accuracy: {test_accuracy:.2f}%"
40
41 # Concatenate all the batches
42 all_targets = np.concatenate(all_targets)
43 all_probs = np.concatenate(all_probs)
44
45 # Compute ROC curve and AUC
46 fpr, tpr, thresholds = roc_curve(all_targets, all_probs)
```
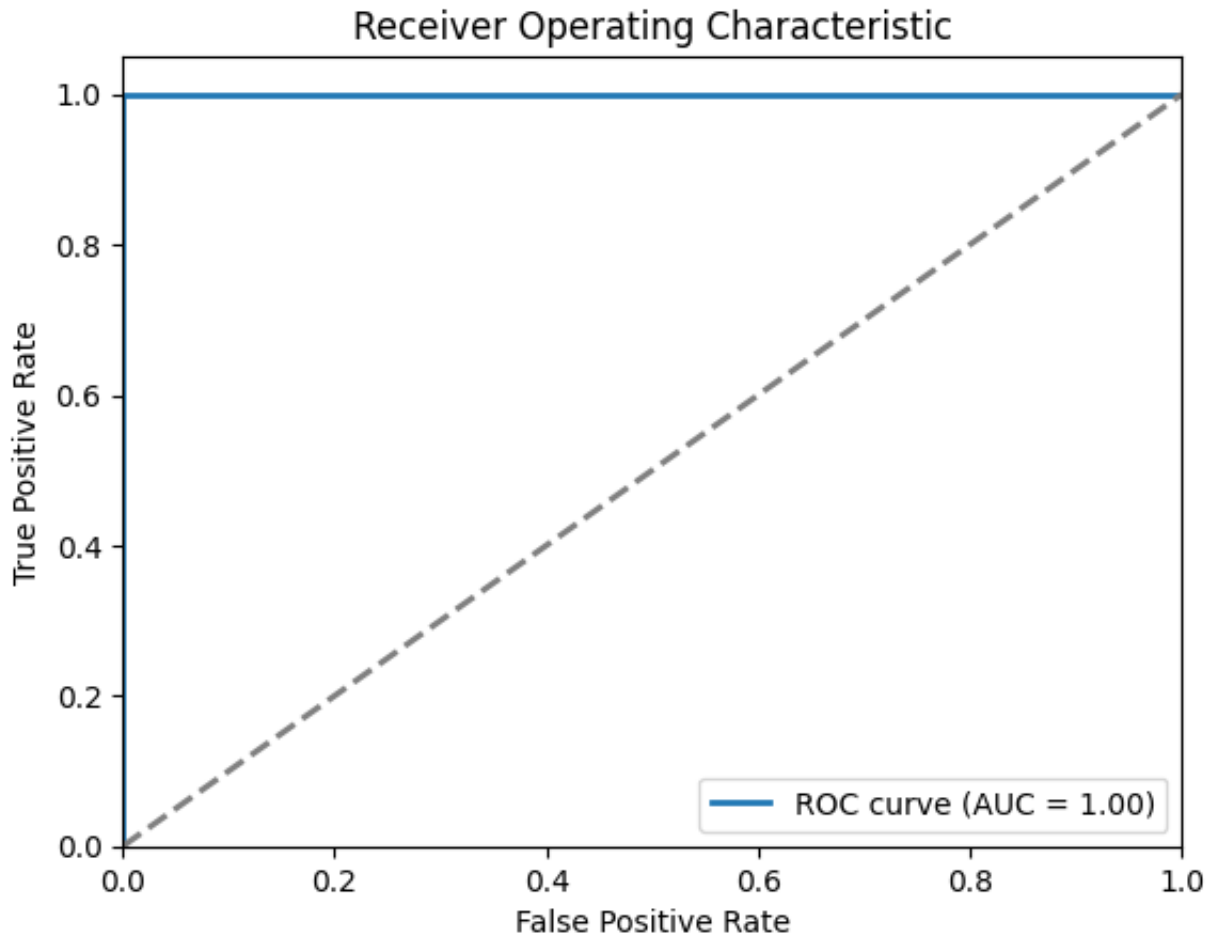
```python
47 roc_auc = auc(fpr, tpr)
48 print("AUC: {:.4f}".format(roc_auc))
49
50 # Plot ROC curve
51 plt.figure()
52 plt.plot(fpr, tpr, lw=2, label='ROC curve (AUC = %0.2f)' % roc_auc)
53 plt.plot([0, 1], [0, 1], lw=2, linestyle='--', color='gray')
54 plt.xlim([0.0, 1.0])
55 plt.ylim([0.0, 1.05])
56 plt.xlabel('False Positive Rate')
57 plt.ylabel('True Positive Rate')
58 plt.title('Receiver Operating Characteristic')
59 plt.legend(loc="lower right")
60 plt.show()
61
62 # Save final model
63 torch.save({
64     'model_state_dict': vit.state_dict(),
65     'test_accuracy': test_accuracy,
66     'test_loss': avg_test_loss
67 }, f"{save_dir}/vit_final_model.pt")
68
69 print("Test completed!")
70
```

```
Evaluating on test set...
Testing: 100%|██████████████| 2/2 [00:00<00:00,  2.61it/s]Test set - Loss: 0.01006
AUC: 1.0000
```



Receiver Operating Characteristic

```
Test completed!
```

```
 1 # Code for visualization of training data, not used in actual training but ke|
 2 non_zero_values = data1[data1 != 0]
 3
 4 # Calculate statistics
 5 min_non_zero = np.min(non_zero_values)
 6 max_non_zero = np.max(non_zero_values)
 7 mean_non_zero = np.mean(non_zero_values)
 8 median_non_zero = np.median(non_zero_values)
 9 std_non_zero = np.std(non_zero_values)
10
11 # Create figure with subplots
12 fig, (ax1, ax2) = plt.subplots(2, 1, figsize=(10, 12))
13
```

```python
14 # Histogram with KDE
15 sns.histplot(non_zero_values, kde=True, ax=ax1)
16 ax1.set_title('Distribution of Non-Zero Values', fontsize=14)
17 ax1.set_xlabel('Value', fontsize=12)
18 ax1.set_ylabel('Frequency', fontsize=12)
19 ax1.axvline(min_non_zero, color='r', linestyle='--', label=f'Min: {min_non_ze
20 ax1.axvline(max_non_zero, color='g', linestyle='--', label=f'Max: {max_non_ze
21 ax1.axvline(mean_non_zero, color='b', linestyle='-', label=f'Mean: {mean_non_
22 ax1.axvline(median_non_zero, color='purple', linestyle='-.', label=f'Median:
23 ax1.legend()
24
25 # Boxplot
26 sns.boxplot(x=non_zero_values, ax=ax2)
27 ax2.set_title('Boxplot of Non-Zero Values', fontsize=14)
28 ax2.set_xlabel('Value', fontsize=12)
29
30 # Add text with statistics
31 stats_text = (f"Non-zero count: {len(non_zero_values)}\n"
32               f"Min: {min_non_zero:.2f}\n"
33               f"Max: {max_non_zero:.2f}\n"
34               f"Mean: {mean_non_zero:.2f}\n"
35               f"Median: {median_non_zero:.2f}\n"
36               f"Std Dev: {std_non_zero:.2f}")
37
38 fig.text(0.15, 0.01, stats_text, fontsize=12, bbox=dict(facecolor='white', al
39
40 plt.tight_layout(rect=[0, 0.03, 1, 0.97])
41 plt.savefig('non_zero_distribution.png')
42 plt.show()
43
44 # Print summary to console
45 print(f"Non-zero value summary:")
46 print(f"Count: {len(non_zero_values)}")
47 print(f"Min: {min_non_zero:.2f}")
48 print(f"Max: {max_non_zero:.2f}")
49 print(f"Mean: {mean_non_zero:.2f}")
50 print(f"Median: {median_non_zero:.2f}")
51 print(f"Standard deviation: {np.std(non_zero_values):.2f}")
```