

```

1 import torch
2 from torch import nn
3
4 from einops import rearrange, repeat
5 from einops.layers.torch import Rearrange
6
7 # helpers
8
9 def pair(t):
10     return t if isinstance(t, tuple) else (t, t)
11
12 # classes
13
14 class PreNorm(nn.Module):
15     def __init__(self, dim, fn):
16         super().__init__()
17         self.norm = nn.LayerNorm(dim)
18         self.fn = fn
19     def forward(self, x, **kwargs):
20         return self.fn(self.norm(x), **kwargs)
21
22 class FeedForward(nn.Module):
23     def __init__(self, dim, hidden_dim, dropout = 0.):
24         super().__init__()
25         self.net = nn.Sequential(
26             nn.Linear(dim, hidden_dim),
27             nn.GELU(),
28             nn.Dropout(dropout),
29             nn.Linear(hidden_dim, dim),
30             nn.Dropout(dropout)
31         )
32     def forward(self, x):
33         return self.net(x)

```

PreNorm : LayerNorm

FeedForward: mlp

```

35 class Attention(nn.Module):
36     def __init__(self, dim, heads = 8, dim_head = 64, dropout = 0.):
37         super().__init__()
38         inner_dim = dim_head * heads
39         project_out = not (heads == 1 and dim_head == dim)
40
41         self.heads = heads
42         self.scale = dim_head ** -0.5
43
44         self.attend = nn.Softmax(dim = -1)
45         self.dropout = nn.Dropout(dropout)
46
47         self.to_qkv = nn.Linear(dim, inner_dim * 3, bias = False)
48
49         self.to_out = nn.Sequential(
50             nn.Linear(inner_dim, dim),
51             nn.Dropout(dropout)
52         ) if project_out else nn.Identity()
53
54     def forward(self, x):
55         qkv = self.to_qkv(x).chunk(3, dim = -1)
56         q, k, v = map(lambda t: rearrange(t, 'b n (h d) -> b h n d', h = self.heads), qkv)
57
58         dots = torch.matmul(q, k.transpose(-1, -2)) * self.scale
59
60         attn = self.attend(dots)
61         attn = self.dropout(attn)
62
63         out = torch.matmul(attn, v)
64         out = rearrange(out, 'b h n d -> b n (h d)')
65         return self.to_out(out)

```

$q, k, v$  linear projection  
 split head  
 $\text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)v$   
 concat head

Attention :

einops optimization

```

66 class Transformer(nn.Module):
67     def __init__(self, dim, depth, heads, dim_head, mlp_dim, dropout = 0.):
68         super().__init__()
69         self.layers = nn.ModuleList([])
70         for _ in range(depth):
71             self.layers.append(nn.ModuleList([
72                 PreNorm(dim, Attention(dim, heads = heads, dim_head = dim_head, dropout = dropout)),
73                 PreNorm(dim, FeedForward(dim, mlp_dim, dropout = dropout))
74             ]))
75
76     def forward(self, x):
77         for attn, ff in self.layers:
78             x = attn(x) + x
79             x = ff(x) + x
80         return x

```

residual connection

Transformer : only encoder  
attention + feed forward

```

82 class ViT(nn.Module):
83     def __init__(self, *, image_size, patch_size, num_classes, dim, depth, heads, mlp_dim, pool = 'cls', channels = 3, dim_head = 64, dropout
84         super().__init__()
85         image_height, image_width = pair(image_size)
86         patch_height, patch_width = pair(patch_size)
87
88         assert image_height % patch_height == 0 and image_width % patch_width == 0, 'Image dimensions must be divisible by the patch size.'
89
90         num_patches = (image_height // patch_height) * (image_width // patch_width)
91         patch_dim = channels * patch_height * patch_width
92         assert pool in {'cls', 'mean'}, 'pool type must be either cls (cls token) or mean (mean pooling)'
93
94         self.to_patch_embedding = nn.Sequential(
95             Rearrange('b c (h p1) (w p2) -> b (h w) (p1 p2 c)', p1 = patch_height, p2 = patch_width), } patchify
96             nn.LayerNorm(patch_dim),
97             nn.Linear(patch_dim, dim),
98             nn.LayerNorm(dim),
99         )
100
101         self.pos_embedding = nn.Parameter(torch.randn(1, num_patches + 1, dim))
102         self.cls_token = nn.Parameter(torch.randn(1, 1, dim)) ← learnable classification token
103         self.dropout = nn.Dropout(emb_dropout)
104
105         self.transformer = Transformer(dim, depth, heads, dim_head, mlp_dim, dropout)
106
107         self.pool = pool
108         self.to_latent = nn.Identity()
109
110         self.mlp_head = nn.Sequential(
111             nn.LayerNorm(dim),
112             nn.Linear(dim, num_classes)
113         )
114
115     def forward(self, img):
116         x = self.to_patch_embedding(img)
117         b, n, _ = x.shape
118
119         cls_tokens = repeat(self.cls_token, '1 1 d -> b 1 d', b = b) } embedding
120         x = torch.cat((cls_tokens, x), dim=1)
121         x += self.pos_embedding[:, :(n + 1)]
122         x = self.dropout(x)
123
124         x = self.transformer(x) ← transformer
125
126         x = x.mean(dim = 1) if self.pool == 'mean' else x[:, 0] ← cls token / mean
127
128         x = self.to_latent(x)
129         return self.mlp_head(x)

```

ViT :

patch\_embedding +  
cls\_token +  
positional\_embedding +  
transformer +  
mlp  
(learnable)