

```

1 class VisionTransformer(nn.Module):
2     def __init__(self, img_size, patch_size, embed_dim, num_heads, num_layers, num_classes):
3         super(VisionTransformer, self).__init__()

4         self.num_patches = (img_size // patch_size) ** 2
5         self.patch_size = patch_size
6         self.patch_dim = 3 * patch_size * patch_size # 3 color channels

7
8         self.patch_embed = nn.Linear(self.patch_dim, embed_dim)
9         self.cls_token = nn.Parameter(torch.randn(1, 1, embed_dim))
10        self.pos_embed = nn.Parameter(torch.randn(1, self.num_patches + 1, embed_dim))

11
12        self.transformer = nn.TransformerEncoder(
13            nn.TransformerEncoderLayer(embed_dim, num_heads, dim_feedforward=embed_dim*4),
14            num_layers
15        )
16        self.mlp_head = nn.Linear(embed_dim, num_classes)

17    def forward(self, x):
18        B, _, _, _ = x.shape # Get batch size. We don't need the channels, height, or width
19        x = x.unfold(2, self.patch_size, self.patch_size).unfold(3, self.patch_size, self.patch_size)
20        x = x.permute(0, 2, 3, 1, 4, 5).contiguous().view(B, self.num_patches, -1)
21        x = self.patch_embed(x)
22
23        cls_tokens = self.cls_token.expand(B, -1, -1)

```

```

1 class VisionTransformer(nn.Module):
2     def __init__(self, img_size, patch_size, embed_dim, num_heads, num_layers, num_classes):
3         super(VisionTransformer, self).__init__()
4
5         self.num_patches = (img_size // patch_size) ** 2
6         self.patch_size = patch_size
7         self.embed_dim = embed_dim
8
9         # Conv2d-based Patch Embedding (Better Stability)
10        self.patch_embed = nn.Conv2d(3, embed_dim, kernel_size=patch_size, stride=patch_size)
11
12        # Class Token & Positional Embedding (Fixed Initialization)
13        self.cls_token = nn.Parameter(torch.zeros(1, 1, embed_dim))
14        self.pos_embed = nn.Parameter(torch.zeros(1, self.num_patches + 1, embed_dim))
15
16        # Transformer with Pre-Normalization (More Stable Training)
17        encoder_layer = nn.TransformerEncoderLayer(embed_dim, num_heads, dim_feedforward=embed_dim*4,
18            dropout=0.1, activation="gelu", batch_first=True)
19        self.transformer = nn.TransformerEncoder(encoder_layer, num_layers)
20
21        # Classification Head (LayerNorm before Output Improves Stability)
22        self.norm = nn.LayerNorm(embed_dim)
23
24        self.mlp_head = nn.Linear(embed_dim, num_classes)
25
26        # Initialize Weights
27        self._init_weights()
28
29        # Initialize Weights with Truncated Normal Distribution
30        def _init_weights(self):
31            trunc_normal_(self.pos_embed, std=0.02)
32            trunc_normal_(self.cls_token, std=0.02)
33            self.apply(self._init_mlp_weights)
34
35        # Initialize Weights for MLP Head
36        def _init_mlp_weights(self, m):
37            if isinstance(m, nn.Linear):
38                trunc_normal_(m.weight, std=0.02)
39            if m.bias is not None:
40                nn.init.constant_(m.bias, 0)
41            elif isinstance(m, nn.LayerNorm):
42                nn.init.constant_(m.bias, 0)
43                nn.init.constant_(m.weight, 1.0)
44
45    def forward(self, x):
46        B = x.shape[0]
47        # Patch Embedding (Conv2d)

```

```

24 x = torch.cat((cls_tokens, x), dim=1)
25 x += self.pos_embed
26
27
27 x = self.transformer(x)
28 x = x[:, 0, :]
29
29 return self.mlp_head(x)
30
31 def build_transformer(config):
32 return VisionTransformer(
33 config['img_size'], config['patch_size'], config['embed_dim'], config['num_heads'],
34 config['num_layers'], config['num_classes']
35 )

```

```

48 x = self.patch_embed(x) # [B, embed_dim, H/patch_size, W/patch_size]
49 x = x.flatten(2).transpose(1, 2) # [B, num_patches, embed_dim]
50
51 # Class Token
52 cls_tokens = self.cls_token.expand(B, -1, -1) # [B, 1, embed_dim]
53 x = torch.cat((cls_tokens, x), dim=1) # [B, num_patches+1, embed_dim]
54
55 # Positional Encoding
56 x = x + self.pos_embed[:, :x.shape[1], :]
57
58 # Transformer Encoder with Pre-Normalization
59 x = self.transformer(x)
60
61 # Classification Head
62 x = self.norm(x[:, 0, :]) # Use CLS Token
63 return self.mlp_head(x)
64
65 def build_transformer(config):
66 return VisionTransformer(
67 config['img_size'], config['patch_size'], config['embed_dim'], config['num_heads'],
68 config['num_layers'], config['num_classes']
69 )

```