

DeiT-S Transfer Learning & Fine-Tuning (PyTorch + `timm`)

1 Setup

```
1 pip install torch torchvision timm==1.*
```

Folder layout (ImageFolder):

```
1 data/  
2   train/  
3     class_a/ img001.jpg ...  
4     class_b/ ...  
5   val/  
6     class_a/ ...  
7     class_b/ ...
```

2 Choices & Defaults

- **Backbone:** `deit_small_patch16_224` (or distilled: `deit_small_distilled_patch16_224`).
- **Input:** 224×224, ImageNet normalization.
- **Optim:** AdamW; **Schedule:** warmup + cosine.
- **Augment:** RandAugment/AutoAugment, Mixup, CutMix, label smoothing.
- **LR scaling:** $\text{base_lr} = 5e-4 \times \frac{\text{batch}}{256}$ (full fine-tune).

3 Minimal Working Script (train + eval)

```
1 import torch, timm  
2 from torch import nn  
3 from torch.utils.data import DataLoader  
4 from torchvision import datasets  
5 from timm.data import create_transform  
6 from timm.loss import LabelSmoothingCrossEntropy,  
   SoftTargetCrossEntropy  
7 from timm.optim import create_optimizer_v2  
8 from timm.scheduler import create_scheduler  
9  
10 # ---- Config ----  
11 data_dir = "data"  
12 model_name = "deit_small_patch16_224" # or "  
   deit_small_distilled_patch16_224"  
13 num_classes = 10 # <-- set to your dataset  
14 epochs = 50  
15 batch_size = 64  
16 img_size = 224
```

```

17 use_mixup      = True
18 device        = "cuda" if torch.cuda.is_available() else "cpu"
19
20 # ---- Transforms ----
21 train_tf = create_transform(
22     input_size=img_size, is_training=True,
23     color_jitter=0.4, auto_augment='rand-m9-mstd0.5-inc1',
24     re_prob=0.25, re_mode='pixel', re_count=1, interpolation='bicubic'
25 )
26 val_tf = create_transform(input_size=img_size, is_training=False,
27     interpolation='bicubic')
28
29 # ---- Datasets / Loaders ----
30 train_ds = datasets.ImageFolder(f"{data_dir}/train", transform=train_tf)
31 val_ds   = datasets.ImageFolder(f"{data_dir}/val",   transform=val_tf)
32 train_loader = DataLoader(train_ds, batch_size=batch_size, shuffle=True,
33     num_workers=8, pin_memory=True)
34 val_loader   = DataLoader(val_ds,   batch_size=batch_size, shuffle=False,
35     num_workers=8, pin_memory=True)
36
37 # ---- Model ----
38 model = timm.create_model(model_name, pretrained=True, num_classes=
39     num_classes).to(device)
40
41 # ---- Mixup/CutMix & Loss ----
42 mixup_fn = None
43 if use_mixup:
44     from timm.data.mixup import Mixup
45     mixup_fn = Mixup(mixup_alpha=0.8, cutmix_alpha=1.0, label_smoothing
46         =0.1, num_classes=num_classes)
47     criterion = SoftTargetCrossEntropy()
48 else:
49     criterion = LabelSmoothingCrossEntropy(smoothing=0.1)
50
51 # ---- Optim / Scheduler ----
52 opt = create_optimizer_v2(model, opt='adamw', weight_decay=0.05, lr=5e
53     -4 * (batch_size/256))
54 sched, _ = create_scheduler({'sched':'cosine', 'epochs':epochs, '
55     warmup_epochs':5, 'cooldown_epochs':0, 'min_lr':1e-6}, opt)
56 scaler = torch.cuda.amp.GradScaler(enabled=(device=="cuda"))
57
58 def train_one_epoch():
59     model.train()
60     total, correct, loss_sum = 0, 0, 0.0
61     for images, targets in train_loader:
62         images, targets = images.to(device), targets.to(device)
63         if mixup_fn: images, targets = mixup_fn(images, targets)
64         opt.zero_grad(set_to_none=True)
65         with torch.cuda.amp.autocast(enabled=(device=="cuda")):
66             outputs = model(images)
67             loss = criterion(outputs, targets)
68             scaler.scale(loss).backward()
69             scaler.step(opt); scaler.update()
70             loss_sum += loss.item() * images.size(0)
71         if not mixup_fn:
72             pred = outputs.argmax(1); correct += (pred == targets).sum
73             ().item()

```

```

66         total += images.size(0)
67     sched.step(None)
68     acc = (correct/total*100) if not mixup_fn else float('nan')
69     return loss_sum/total, acc
70
71 @torch.no_grad()
72 def evaluate():
73     model.eval()
74     total, correct, loss_sum = 0, 0, 0.0
75     for images, targets in val_loader:
76         images, targets = images.to(device), targets.to(device)
77         with torch.cuda.amp.autocast(enabled=(device=="cuda")):
78             outputs = model(images)
79             loss = nn.CrossEntropyLoss()(outputs, targets)
80             loss_sum += loss.item() * images.size(0)
81             pred = outputs.argmax(1); correct += (pred == targets).sum().item()
82             total += images.size(0)
83     return loss_sum/total, correct/total*100
84
85 best_acc, best_path = 0.0, "best_deit_s.pth"
86 for epoch in range(1, epochs+1):
87     tr_loss, tr_acc = train_one_epoch()
88     val_loss, val_acc = evaluate()
89     if val_acc > best_acc:
90         best_acc = val_acc
91         torch.save({'model': model.state_dict()}, best_path)
92     print(f"Epoch {epoch:03d} | train {tr_loss:.4f}/{tr_acc:.2f} | val {val_loss:.4f}/{val_acc:.2f}")

```

4 Strategy A: Feature Extraction (freeze backbone)

Use when data are scarce / risk of overfitting is high.

```

1  # Freeze everything except classifier head(s)
2  for n, p in model.named_parameters():
3      if not (n.startswith('head') or n.startswith('fc') or 'distill' in n):
4          p.requires_grad = False
5
6  # Train small head with a slightly higher LR and no weight decay
7  opt = create_optimizer_v2(
8      filter(lambda p: p.requires_grad, model.parameters()),
9      opt='adamw', weight_decay=0.0, lr=1e-3
10 )

```

Progressive unfreezing: after 5–10 epochs, unfreeze the last transformer block; then more blocks if validation improves.

5 Strategy B: Full Fine-Tuning (+ Layer-Wise LR Decay)

```

1  def param_groups_llrd(model, base_lr=5e-4, weight_decay=0.05,
2      decay_rate=0.75):
3      layers = []
4      layers.append([*model.patch_embed.parameters(), model.pos_embed])
5      # shallowest

```

```

4     for blk in model.blocks: layers.append(list(blk.parameters()))
5     layers.append(list(model.head.parameters()))
6     # head
7     n = len(layers); groups = []
8     for i, params in enumerate(layers):
9         lr = base_lr * (decay_rate ** (n - i - 1))
10        groups.append({'params': params, 'lr': lr, 'weight_decay':
11                        weight_decay})
12    return groups
13
14 opt = torch.optim.AdamW(param_groups_llrd(model), betas=(0.9, 0.999))

```

6 Distilled Variant (Optional)

If you choose `deit_small_distilled_patch16_224`, a distillation token/head is present. With pretrained distilled weights you can fine-tune *without* a teacher; train as usual (timm handles the two heads).

7 Hyperparameters That Usually Work

- Epochs: 50–100 (30–50 for smaller datasets).
- Batch size: as large as fits (32–256). Scale LR accordingly.
- LR: $5e-4$ full fine-tune; $1e-3$ head-only; 5 warmup epochs; cosine decay.
- Weight decay: 0.05 (head-only: 0.0–0.02).
- Augment: RandAugment; Mixup=0.8; CutMix=1.0; label smoothing=0.1.
- Regularization: keep DeiT-S default `drop_path` (~ 0.1); use AMP.
- Early stopping: monitor val loss/acc (patience ≈ 10).

8 Evaluation Tips

Track top-1/top-5, confusion matrix, macro-F1 for imbalance. Consider EMA of weights (`ModelEmaV2` in `timm`) for stability.

9 Export / Deployment

```

1 model.eval()
2 example = torch.randn(1, 3, 224, 224).to(device)
3 traced = torch.jit.trace(model, example)
4 traced.save("deit_s_traced.pt")

```

10 Sanity Checklist

- ☐ Correct `num_classes`.
- ☐ Proper `data/` layout and splits.

- ☐ 224 input, ImageNet normalization.
- ☐ LR scaled to batch; warmup on; cosine decay.
- ☐ AMP enabled; best checkpoint saved.
- ☐ Final test on a held-out set.