

## Lab Session 4.Vision-Language Model (VLM) Prompt Tuning

### Contents

- [1] Preparation
- [2] Load pre-trained CLIP Model
- [3] CoOpCLIP Implementation
- [4] CoOpCLIP Training

### References

- Learning to Prompt for Vision-Language Models (CoOp): <https://github.com/KaiyangZhou/CoOp>
- Prompt Learning via Meta-Regularization (ProMetaR): <https://github.com/mlvlab/ProMetaR>

### ✓ [1] Preparation

#### ✓ 1. Clone github repository

```
1 !git clone https://github.com/mlvlab/ProMetaR.git
```

 fatal: destination path 'ProMetaR' already exists and is not an empty directory.

- It will make ProMetaR folder on left side.

#### ✓ 2. Install Requirements

```
1 %cd ProMetaR/  
2  
3 !git clone https://github.com/KaiyangZhou/Dassl.pytorch.git  
4 %cd Dassl.pytorch/  
5  
6 # Install dependencies  
7 !pip install -r requirements.txt  
8 !cp -r dassl ../  
9 # Install this library (no need to re-build if the source code is modified)  
10 # !python setup.py develop  
11 %cd ..  
12  
13 !pip install -r requirements.txt
```



```

Requirement already satisfied: certifi<2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests->learn2learn==0.2.0->r requirements.tx
Requirement already satisfied: ecos>=2 in /usr/local/lib/python3.10/dist-packages (from cvxpy>=1.1.0->qpth>=0.0.15->learn2learn==0.2.0->r requ
Requirement already satisfied: clarabel>=0.5.0 in /usr/local/lib/python3.10/dist-packages (from cvxpy>=1.1.0->qpth>=0.0.15->learn2learn==0.2.0->r
Requirement already satisfied: scs>=3.2.4.post1 in /usr/local/lib/python3.10/dist-packages (from cvxpy>=1.1.0->qpth>=0.0.15->learn2learn==0.2.0->r
Requirement already satisfied: boto>=2.29.1 in /usr/local/lib/python3.10/dist-packages (from gcs-oauth2-boto-plugin>=3.2->gsutil->learn2learn==0.2.
Requirement already satisfied: oauth2client>=2.2.0 in /usr/local/lib/python3.10/dist-packages (from gcs-oauth2-boto-plugin>=3.2->gsutil->learn2lear
Requirement already satisfied: pyasn1>=0.1.3 in /usr/local/lib/python3.10/dist-packages (from rsa<5,>=3.1.4->google-auth==2.17.0->google-auth[aiht
Requirement already satisfied: pyu2f in /usr/local/lib/python3.10/dist-packages (from google-reauth>=0.1.0->gsutil->learn2learn==0.2.0->-r requir
Requirement already satisfied: cryptography<44,>=41.0.5 in /usr/local/lib/python3.10/dist-packages (from pyOpenSSL>=0.13->gsutil->learn2learn==0.2.
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/dist-packages (from jinja2->torch>=1.1.0->learn2learn==0.2.0->-r requir
Requirement already satisfied: aiohappyeyeballs>=2.3.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp<4.0.0dev,>=3.6.2->google-auth[aiht
Requirement already satisfied: aiosignal>=1.1.2 in /usr/local/lib/python3.10/dist-packages (from aiohttp<4.0.0dev,>=3.6.2->google-auth[aihttp]=2.
Requirement already satisfied: attrs>=17.3.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp<4.0.0dev,>=3.6.2->google-auth[aihttp]=2.17.
Requirement already satisfied: frozenlist>=1.1.1 in /usr/local/lib/python3.10/dist-packages (from aiohttp<4.0.0dev,>=3.6.2->google-auth[aihttp]=2.
Requirement already satisfied: multidict<7.0,>=4.5 in /usr/local/lib/python3.10/dist-packages (from aiohttp<4.0.0dev,>=3.6.2->google-auth[aihttp]=
Requirement already satisfied: propcache>=0.2.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp<4.0.0dev,>=3.6.2->google-auth[aihttp]=2.
Requirement already satisfied: yarl<2.0,>=1.17.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp<4.0.0dev,>=3.6.2->google-auth[aihttp]=2.
Requirement already satisfied: async-timeout<6.0,>=4.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp<4.0.0dev,>=3.6.2->google-auth[aiht
Requirement already satisfied: cffi>=1.12 in /usr/local/lib/python3.10/dist-packages (from cryptography<44,>=41.0.5->pyOpenSSL>=0.13->gsutil->learr
Requirement already satisfied: qdldl in /usr/local/lib/python3.10/dist-packages (from osqp>=0.6.2->cvxpy>=1.1.0->qpth>=0.0.15->learn2learn==0.2.0->
Requirement already satisfied: pycparser in /usr/local/lib/python3.10/dist-packages (from cffi>=1.12->cryptography<44,>=41.0.5->pyOpenSSL>=0.13->gs

```

- If an error occurs, click 'Run Session Again' and then restart the runtime from the beginning.

### 3. Load Requirements and functions

```

1 import os.path as osp
2 from collections import OrderedDict
3 import math
4 import torch
5 import torch.nn as nn
6 from torch.nn import functional as F
7 from torch.cuda.amp import GradScaler, autocast
8 from PIL import Image
9 import torchvision.transforms as transforms
10 import torch
11 from clip import clip
12 from clip.simple_tokenizer import SimpleTokenizer as _Tokenizer
13 import time
14 from tqdm import tqdm
15 import datetime
16 import argparse
17 from dassl.utils import setup_logger, set_random_seed, collect_env_info
18 from dassl.config import get_cfg_default
19 from dassl.engine import build_trainer
20 from dassl.engine import TRAINER_REGISTRY, TrainerX
21 from dassl.metrics import compute_accuracy
22 from dassl.utils import load_pretrained_weights, load_checkpoint
23 from dassl.optim import build_optimizer, build_lr_scheduler
24
25 # custom
26 import datasets.oxford_pets
27 import datasets.oxford_flowers
28 import datasets.fgvc_aircraft
29 import datasets.dtd
30 import datasets.eurosat
31 import datasets.stanford_cars
32 import datasets.food101
33 import datasets.sun397
34 import datasets.caltech101
35 import datasets.ucf101
36 import datasets.imagenet
37 import datasets.imagenet_sketch
38 import datasets.imagenetv2
39 import datasets.imagenet_a
40 import datasets.imagenet_r

```

```

1 def print_args(args, cfg):
2     print("*****")
3     print("** Arguments **")
4     print("*****")
5     optkeys = list(args.__dict__.keys())
6     optkeys.sort()
7     for key in optkeys:
8         print("{}: {}".format(key, args.__dict__[key]))
9     print("*****")
10    print("** Config **")
11    print("*****")
12    print(cfg)
13
14 def reset_cfg(cfg, args):

```

```

15     if args.root:
16         cfg.DATASET.ROOT = args.root
17     if args.output_dir:
18         cfg.OUTPUT_DIR = args.output_dir
19     if args.seed:
20         cfg.SEED = args.seed
21     if args.trainer:
22         cfg.TRAINER.NAME = args.trainer
23     cfg.DATASET.NUM_SHOTS = 16
24     cfg.DATASET.SUBSAMPLE_CLASSES = args.subsample_classes
25     cfg.DATALOADER.TRAIN_X.BATCH_SIZE = args.train_batch_size
26     cfg.OPTIM.MAX_EPOCH = args.epoch
27
28 def extend_cfg(cfg):
29     """
30     Add new config variables.
31     """
32     from yacs.config import CfgNode as CN
33     cfg.TRAINER.COOP = CN()
34     cfg.TRAINER.COOP.N_CTX = 16 # number of context vectors
35     cfg.TRAINER.COOP.CSC = False # class-specific context
36     cfg.TRAINER.COOP.CTX_INIT = "" # initialization words
37     cfg.TRAINER.COOP.PREC = "fp16" # fp16, fp32, amp
38     cfg.TRAINER.COOP.CLASS_TOKEN_POSITION = "end" # 'middle' or 'end' or 'front'
39     cfg.TRAINER.COOCOOP = CN()
40     cfg.TRAINER.COOCOOP.N_CTX = 4 # number of context vectors
41     cfg.TRAINER.COOCOOP.CTX_INIT = "a photo of a" # initialization words
42     cfg.TRAINER.COOCOOP.PREC = "fp16" # fp16, fp32, amp
43     cfg.TRAINER.PROMETAR = CN()
44     cfg.TRAINER.PROMETAR.N_CTX_VISION = 4 # number of context vectors at the vision branch
45     cfg.TRAINER.PROMETAR.N_CTX_TEXT = 4 # number of context vectors at the language branch
46     cfg.TRAINER.PROMETAR.CTX_INIT = "a photo of a" # initialization words
47     cfg.TRAINER.PROMETAR.PREC = "fp16" # fp16, fp32, amp
48     cfg.TRAINER.PROMETAR.PROMPT_DEPTH_VISION = 9 # Max 12, minimum 0, for 0 it will be using shallow IVLP prompting (J=1)
49     cfg.TRAINER.PROMETAR.PROMPT_DEPTH_TEXT = 9 # Max 12, minimum 0, for 0 it will be using shallow IVLP prompting (J=1)
50     cfg.DATASET.SUBSAMPLE_CLASSES = "all" # all, base or new
51     cfg.TRAINER.PROMETAR.ADAPT_LR = 0.0005
52     cfg.TRAINER.PROMETAR.LR_RATIO = 0.0005
53     cfg.TRAINER.PROMETAR.FAST_ADAPTATION = False
54     cfg.TRAINER.PROMETAR.MIXUP_ALPHA = 0.5
55     cfg.TRAINER.PROMETAR.MIXUP_BETA = 0.5
56     cfg.TRAINER.PROMETAR.DIM_RATE=8
57     cfg.OPTIM_VNET = CN()
58     cfg.OPTIM_VNET.NAME = "adam"
59     cfg.OPTIM_VNET.LR = 0.0003
60     cfg.OPTIM_VNET.WEIGHT_DECAY = 5e-4
61     cfg.OPTIM_VNET.MOMENTUM = 0.9
62     cfg.OPTIM_VNET.SGD_DAMPNING = 0
63     cfg.OPTIM_VNET.SGD_NESTEROV = False
64     cfg.OPTIM_VNET.RMSPROP_ALPHA = 0.99
65     cfg.OPTIM_VNET.ADM_BETA1 = 0.9
66     cfg.OPTIM_VNET.ADM_BETA2 = 0.999
67     cfg.OPTIM_VNET.STAGED_LR = False
68     cfg.OPTIM_VNET.NEW_LAYERS = ()
69     cfg.OPTIM_VNET.BASE_LR_MULT = 0.1
70     # Learning rate scheduler
71     cfg.OPTIM_VNET.LR_SCHEDULER = "single_step"
72     # -1 or 0 means the stepsize is equal to max_epoch
73     cfg.OPTIM_VNET.STEPSIZE = (-1, )
74     cfg.OPTIM_VNET.GAMMA = 0.1
75     cfg.OPTIM_VNET.MAX_EPOCH = 10
76     # Set WARMUP_EPOCH larger than 0 to activate warmup training
77     cfg.OPTIM_VNET.WARMUP_EPOCH = -1
78     # Either linear or constant
79     cfg.OPTIM_VNET.WARMUP_TYPE = "linear"
80     # Constant learning rate when type=constant
81     cfg.OPTIM_VNET.WARMUP_CONS_LR = 1e-5
82     # Minimum learning rate when type=linear
83     cfg.OPTIM_VNET.WARMUP_MIN_LR = 1e-5
84     # Recount epoch for the next scheduler (last_epoch=-1)
85     # Otherwise last_epoch=warmup_epoch
86     cfg.OPTIM_VNET.WARMUP_RECOUNT = True
87
88 def setup_cfg(args):
89     cfg = get_cfg_default()
90     extend_cfg(cfg)
91     # 1. From the dataset config file
92     if args.dataset_config_file:
93         cfg.merge_from_file(args.dataset_config_file)
94     # 2. From the method config file
95     if args.config_file:
96         cfg.merge_from_file(args.config_file)
97     # 3. From input arguments
98     reset_cfg(cfg, args)

```

```

99     cfg.freeze()
100    return cfg

```

### 3. Download EuroSAT dataset

```

1 %mkdir outputs
2 %mkdir data
3
4 %cd data
5 %mkdir eurosat
6 !wget http://madm.dfki.de/files/sentinel/EuroSAT.zip EuroSAT.zip
7
8 !unzip -o EuroSAT.zip -d eurosat/
9 %cd eurosat
10 !gdown 1Ip7yaCWF10ea0FUGga0IUdVi_DDQth1o
11
12 %cd ../..
13

```

```

↻ inflating: eurosat/2750/PermanentCrop/PermanentCrop_615.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1398.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_163.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_970.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_502.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_2472.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1567.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1915.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_2013.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_828.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1106.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1670.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1211.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_2304.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_273.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1088.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_612.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1438.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_164.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1059.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_505.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_977.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_2475.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1912.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1560.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_2014.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1101.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1677.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_19.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1216.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_2303.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1753.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1332.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1495.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_2227.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_118.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1444.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1836.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_2130.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1782.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_579.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1025.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_2409.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_853.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_421.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_386.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_2068.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_882.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_357.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_1.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_65.jpg
inflating: eurosat/2750/PermanentCrop/PermanentCrop_736.jpg
/content/ProMetaR/data/eurosat
Downloading...
From: https://drive.google.com/uc?id=1Ip7yaCWF10ea0FUGga0IUdVi\_DDQth1o
To: /content/ProMetaR/data/eurosat/split_zhou_EuroSAT.json
100% 3.01M/3.01M [00:00<00:00, 23.4MB/s]
/content/ProMetaR

```

The structure of the data folder inside the ProMetaR folder should be as follows:

- eurosat/2750
- eurosat/split\_zhou\_EuroSAT.json

- ✓ [2] Load pre-trained CLIP Model

```

1 _tokenizer = _Tokenizer()
2
3 def load_clip_to_cpu(cfg): # Load CLIP
4     backbone_name = cfg.MODEL.BACKBONE.NAME
5     url = clip._MODELS[backbone_name]
6     model_path = clip._download(url)
7
8     try:
9         # loading JIT archive
10        model = torch.jit.load(model_path, map_location="cpu").eval()
11        state_dict = None
12
13    except RuntimeError:
14        state_dict = torch.load(model_path, map_location="cpu")
15
16    if cfg.TRAINER.NAME == "":
17        design_trainer = "CoOp"
18    else:
19        design_trainer = cfg.TRAINER.NAME
20    design_details = {"trainer": design_trainer,
21                     "vision_depth": 0,
22                     "language_depth": 0, "vision_ctx": 0,
23                     "language_ctx": 0}
24    model = clip.build_model(state_dict or model.state_dict(), design_details)
25
26    return model

```

100% | 351M/351M [00:07<00:00, 44.4MiB/s]

- ✓ Check CLIP model

```
1 print(clip_model)
```



```

        (attn): MultiheadAttention
        (out_proj): NonDynamicallyQuantizableLinear(in_features=512, out_features=512, bias=True)
    )
    (ln_1): LayerNorm((512,), eps=1e-05, elementwise_affine=True)
    (mlp): Sequential(
      (c_fc): Linear(in_features=512, out_features=2048, bias=True)
      (gelu): QuickGELU()
      (c_proj): Linear(in_features=2048, out_features=512, bias=True)
    )
    (ln_2): LayerNorm((512,), eps=1e-05, elementwise_affine=True)
  )
)
)
(token_embedding): Embedding(49408, 512)
(ln_final): LayerNorm((512,), eps=1e-05, elementwise_affine=True)
)

```

### ✓ [3] CoOpCLIP Implementation

CoOpCLIP is composed of pre-trained CLIP Text encoder, pre-trained CLIP Image encoder, and learnable prompt.

#### ✓ Make Module 1: CLIP Text Encoder

##### Input

- token prefix (SOS token) + learnable prompt + class label + token suffix (CLS token)

##### Output

- text feature of input prompts

```

1 class TextEncoder(nn.Module):
2     def __init__(self, clip_model):
3         super().__init__()
4         self.transformer = clip_model.transformer
5         self.positional_embedding = clip_model.positional_embedding
6         self.ln_final = clip_model.ln_final
7         self.text_projection = clip_model.text_projection
8         self.dtype = clip_model.dtype
9
10    def forward(self, prompts, tokenized_prompts): # Call model forward
11        x = prompts + self.positional_embedding.type(self.dtype)
12        x = x.permute(1, 0, 2) # NLD -> LND
13        x = self.transformer(x)
14        x = x.permute(1, 0, 2) # LND -> NLD
15        x = self.ln_final(x).type(self.dtype)
16        x = x[torch.arange(x.shape[0]), tokenized_prompts.argmax(dim=-1)] @ self.text_projection
17        return x

```

#### ✓ Make Module 2: CLIP Image Encoder

##### Input

- image

##### Output

- image feature

```
1 print(clip_model.visual)
```



```

        (mlp): Sequential(
          (c_fc): Linear(in_features=768, out_features=3072, bias=True)
          (gelu): QuickGELU()
          (c_proj): Linear(in_features=3072, out_features=768, bias=True)
        )
        (ln_2): LayerNorm((768.), eps=1e-05, elementwise_affine=True)
      )
    (10): ResidualAttentionBlock(
      (attn): MultiheadAttention(
        (out_proj): NonDynamicallyQuantizableLinear(in_features=768, out_features=768, bias=True)
      )
      (ln_1): LayerNorm((768.), eps=1e-05, elementwise_affine=True)
      (mlp): Sequential(
        (c_fc): Linear(in_features=768, out_features=3072, bias=True)
        (gelu): QuickGELU()
        (c_proj): Linear(in_features=3072, out_features=768, bias=True)
      )
      (ln_2): LayerNorm((768.), eps=1e-05, elementwise_affine=True)
    )
  )
  (ln_post): LayerNorm((768.), eps=1e-05, elementwise_affine=True)
)

```

## ✓ Make Module 3: Learnable Prompt

### Output

#### • Learnable prompt

```

1 class CoOpPromptLearner(nn.Module):
2     def __init__(self, cfg, classnames, clip_model):
3         super().__init__()
4         n_cls = len(classnames)
5         n_ctx = cfg.TRAINER.COOP.N_CTX
6         ctx_init = cfg.TRAINER.COOP.CTX_INIT
7         dtype = clip_model.dtype
8         ctx_dim = clip_model.ln_final.weight.shape[0]
9         clip_imsize = clip_model.visual.input_resolution
10        cfg_imsize = cfg.INPUT.SIZE[0]
11        assert cfg_imsize == clip_imsize, f"cfg_imsize ({cfg_imsize}) must equal to clip_imsize ({clip_imsize})"
12
13        ### Learnable Prompts Initialization ###
14        if ctx_init:
15            # use given words to initialize context vectors
16            ctx_init = ctx_init.replace("_", " ")
17            n_ctx = len(ctx_init.split(" "))
18            prompt = clip.tokenize(ctx_init)
19            with torch.no_grad():
20                embedding = clip_model.token_embedding(prompt).type(dtype)
21            ctx_vectors = embedding[0, 1 : 1 + n_ctx, :]
22            prompt_prefix = ctx_init
23        else:
24            # random initialization
25            if cfg.TRAINER.COOP.CSC:
26                print("Initializing class-specific contexts")
27                ctx_vectors = torch.empty(n_cls, n_ctx, ctx_dim, dtype=dtype)
28            else:
29                print("Initializing a generic context")
30                ctx_vectors = torch.empty(n_ctx, ctx_dim, dtype=dtype)
31            nn.init.normal_(ctx_vectors, std=0.02)
32            prompt_prefix = " ".join(["X"] * n_ctx)
33        print(f'Initial context: "{prompt_prefix}"')
34        print(f"Number of context words (tokens): {n_ctx}")
35        self.ctx = nn.Parameter(ctx_vectors) # Wrap the initialized prompts above as parameters to make them trainable.
36
37        ### Tokenize ###
38        classnames = [name.replace("_", " ") for name in classnames] # 예) "Forest"
39        name_lens = [len(_tokenizer.encode(name)) for name in classnames]
40        prompts = [prompt_prefix + " " + name + "." for name in classnames] # 예) "A photo of Forest."
41        tokenized_prompts = torch.cat([clip.tokenize(p) for p in prompts]) # 예) [49406, 320, 1125, 539...]
42        #####
43
44        with torch.no_grad():
45            embedding = clip_model.token_embedding(tokenized_prompts).type(dtype)

```

```

46     # These token vectors will be saved when in save_model().
47     # but they should be ignored in load_model() as we want to use
48     # those computed using the current class names
49     self.register_buffer("token_prefix", embedding[:, :1, :]) # SOS (문장의 시작을 알려주는 토큰)
50     self.register_buffer("token_suffix", embedding[:, 1 + n_ctx :, :]) # CLS, EOS (문장의 끝을 알려주는 토큰)
51     self.n_cls = n_cls
52     self.n_ctx = n_ctx
53     self.tokenized_prompts = tokenized_prompts # torch.Tensor
54     self.name_lens = name_lens
55
56     def construct_prompts(self, ctx, prefix, suffix, label=None):
57         # dim0 is either batch_size (during training) or n_cls (during testing)
58         # ctx: context tokens, with shape of (dim0, n_ctx, ctx_dim)
59         # prefix: the sos token, with shape of (n_cls, 1, ctx_dim)
60         # suffix: remaining tokens, with shape of (n_cls, *, ctx_dim)
61         if label is not None:
62             prefix = prefix[label]
63             suffix = suffix[label]
64         prompts = torch.cat(
65             [
66                 prefix, # (dim0, 1, dim)
67                 ctx, # (dim0, n_ctx, dim)
68                 suffix, # (dim0, *, dim)
69             ],
70             dim=1,
71         )
72         return prompts
73
74     def forward(self):
75         ctx = self.ctx
76         if ctx.dim() == 2:
77             ctx = ctx.unsqueeze(0).expand(self.n_cls, -1, -1)
78         prefix = self.token_prefix
79         suffix = self.token_suffix
80         prompts = self.construct_prompts(ctx, prefix, suffix) #[시작토큰, Input prompts, 끝 토큰]
81         return prompts

```

## ▼ Make CoOpCLIP (Module1 + Module2 + Module3)

### Input

- Image

### Output

- Logit

### How to compute logit?

- image\_features : Image representation  $\mathbf{f}$
- text\_features : Text representation  $g(\mathbf{t}_i)$
- Logit:  $p(y = i | \mathbf{x}) = \frac{\exp(\cos(g(\mathbf{t}_i), \mathbf{f})/\tau)}{\sum_{j=1}^K \exp(\cos(g(\mathbf{t}_j), \mathbf{f})/\tau)}$

```

1 class CoOpCustomCLIP(nn.Module):
2     def __init__(self, cfg, classnames, clip_model):
3         super().__init__()
4         self.prompt_learner = CoOpPromptLearner(cfg, classnames, clip_model)
5         self.tokenized_prompts = self.prompt_learner.tokenized_prompts
6         self.image_encoder = clip_model.visual
7         self.text_encoder = TextEncoder(clip_model)
8         self.logit_scale = clip_model.logit_scale
9         self.dtype = clip_model.dtype
10
11     def forward(self, image):
12         image_features = self.image_encoder(image.type(self.dtype))
13
14         prompts = self.prompt_learner()
15         tokenized_prompts = self.tokenized_prompts
16         text_features = self.text_encoder(prompts, tokenized_prompts)
17
18         image_features = image_features / image_features.norm(dim=-1, keepdim=True)
19         text_features = text_features / text_features.norm(dim=-1, keepdim=True)
20
21         logit_scale = self.logit_scale.exp()
22         logits = logit_scale * image_features @ text_features.t()
23
24         return logits

```

## ▼ [4] CoOpCLIP Training



## Training configurations

```

1 parser = argparse.ArgumentParser()
2 parser.add_argument("--root", type=str, default="data/", help="path to dataset")
3 parser.add_argument("--output-dir", type=str, default="outputs/cocoop3", help="output directory")
4 parser.add_argument(
5     "--seed", type=int, default=1, help="only positive value enables a fixed seed"
6 )
7 parser.add_argument(
8     "--config-file", type=str, default="configs/trainers/ProMetaR/vit_b16_c2_ep10_batch4_4+ctx.yaml", help="path to config file"
9 )
10 parser.add_argument(
11     "--dataset-config-file",
12     type=str,
13     default="configs/datasets/eurosat.yaml",
14     help="path to config file for dataset setup",
15 )
16 parser.add_argument("--trainer", type=str, default="CoOp", help="name of trainer")
17 parser.add_argument("--eval-only", action="store_true", help="evaluation only")
18 parser.add_argument(
19     "--model-dir",
20     type=str,
21     default="",
22     help="load model from this directory for eval-only mode",
23 )
24 parser.add_argument("--train-batch-size", type=int, default=4)
25 parser.add_argument("--epoch", type=int, default=10)
26 parser.add_argument("--subsample-classes", type=str, default="base")
27 parser.add_argument(
28     "--load-epoch", type=int, default=0, help="load model weights at this epoch for evaluation"
29 )
30 args = parser.parse_args([])

```

## Trainer Class

```

1 @TRAINER_REGISTRY.register(force=True)
2 class CoOp(TrainerX):
3     """Context Optimization (CoOp).
4
5     Learning to Prompt for Vision-Language Models
6     https://arxiv.org/abs/2109.01134
7     """
8
9     def check_cfg(self, cfg):
10         assert cfg.TRAINER.COOP.PREC in ["fp16", "fp32", "amp"]
11
12     def build_model(self):
13         cfg = self.cfg
14         classnames = self.dm.dataset.classnames
15
16         print(f>Loading CLIP (backbone: {cfg.MODEL.BACKBONE.NAME})")
17         clip_model = load_clip_to_cpu(cfg)
18
19         if cfg.TRAINER.COOP.PREC == "fp32" or cfg.TRAINER.COOP.PREC == "amp":
20             # CLIP's default precision is fp16
21             clip_model.float()
22
23         print("Building custom CLIP")
24         self.model = CoOpCustomCLIP(cfg, classnames, clip_model)
25
26         print("Turning off gradients in both the image and the text encoder")
27         for name, param in self.model.named_parameters():
28             if "prompt_learner" not in name:
29                 param.requires_grad_(False)
30
31         if cfg.MODEL.INIT_WEIGHTS:
32             load_pretrained_weights(self.model.prompt_learner, cfg.MODEL.INIT_WEIGHTS)
33
34         self.model.to(self.device)
35         # NOTE: only give prompt_learner to the optimizer
36         self.optim = build_optimizer(self.model.prompt_learner, cfg.OPTIM)
37         self.sched = build_lr_scheduler(self.optim, cfg.OPTIM)
38         self.register_model("prompt_learner", self.model.prompt_learner, self.optim, self.sched)
39
40         self.scaler = GradScaler() if cfg.TRAINER.COOP.PREC == "amp" else None
41
42         # Note that multi-gpu training could be slow because CLIP's size is
43         # big, which slows down the copy operation in DataParallel
44         device_count = torch.cuda.device_count()
45         if device_count > 1:
46             print(f"Multiple GPUs detected (n_gpus={device_count}), use all of them!")

```

```

47         self.model = nn.DataParallel(self.model)
48
49     def before_train(self):
50         directory = self.cfg.OUTPUT_DIR
51         if self.cfg.RESUME:
52             directory = self.cfg.RESUME
53         self.start_epoch = self.resume_model_if_exist(directory)
54
55         # Remember the starting time (for computing the elapsed time)
56         self.time_start = time.time()
57
58     def forward_backward(self, batch):
59         image, label = self.parse_batch_train(batch)
60
61         prec = self.cfg.TRAINER.COOP.PREC
62         output = self.model(image)      # Input image 모델 통과
63         loss = F.cross_entropy(output, label) # Loss 선언
64         self.model_backward_and_update(loss) # Backward 및 모델 parameter 업데이트
65
66         loss_summary = {
67             "loss": loss.item(),
68             "acc": compute_accuracy(output, label)[0].item(),
69         }
70
71         if (self.batch_idx + 1) == self.num_batches:
72             self.update_lr()
73
74         return loss_summary
75
76     def parse_batch_train(self, batch):
77         input = batch["img"]
78         label = batch["label"]
79         input = input.to(self.device)
80         label = label.to(self.device)
81         return input, label
82
83     def load_model(self, directory, epoch=None):
84         if not directory:
85             print("Note that load_model() is skipped as no pretrained model is given")
86             return
87
88         names = self.get_model_names()
89
90         # By default, the best model is loaded
91         model_file = "model-best.pth.tar"
92
93         if epoch is not None:
94             model_file = "model.pth.tar-" + str(epoch)
95
96         for name in names:
97             model_path = osp.join(directory, name, model_file)
98
99             if not osp.exists(model_path):
100                 raise FileNotFoundError('Model not found at "{}"'.format(model_path))
101
102             checkpoint = load_checkpoint(model_path)
103             state_dict = checkpoint["state_dict"]
104             epoch = checkpoint["epoch"]
105
106             # Ignore fixed token vectors
107             if "token_prefix" in state_dict:
108                 del state_dict["token_prefix"]
109
110             if "token_suffix" in state_dict:
111                 del state_dict["token_suffix"]
112
113             print("Loading weights to {} " 'from "{}' (epoch = {})'.format(name, model_path, epoch))
114             # set strict=False
115             self._models[name].load_state_dict(state_dict, strict=False)
116
117     def after_train(self):
118         print("Finish training")
119
120     do_test = not self.cfg.TEST.NO_TEST
121     if do_test:
122         if self.cfg.TEST.FINAL_MODEL == "best_val":
123             print("Deploy the model with the best val performance")
124             self.load_model(self.output_dir)
125         else:
126             print("Deploy the last-epoch model")
127         acc = self.test()
128
129     # Show elapsed time
130     elapsed = round(time.time() - self.time_start)
131     elapsed = str(datetime.timedelta(seconds=elapsed))

```

```

132     print(f"Elapsed: {elapsed}")
133
134     # Close writer
135     self.close_writer()
136     return acc
137
138     def train(self):
139         """Generic training loops."""
140         self.before_train()
141         for self.epoch in range(self.start_epoch, self.max_epoch):
142             self.before_epoch()
143             self.run_epoch()
144             self.after_epoch()
145         acc = self.after_train()
146         return acc
147

```

```

1 def main(args):
2     cfg = setup_cfg(args)
3     if cfg.SEED >= 0:
4         set_random_seed(cfg.SEED)
5
6     if torch.cuda.is_available() and cfg.USE_CUDA:
7         torch.backends.cudnn.benchmark = True
8
9     trainer = build_trainer(cfg)
10    if args.eval_only:
11        trainer.load_model(args.model_dir, epoch=args.load_epoch)
12        acc = trainer.test()
13        return acc
14
15    acc = trainer.train()
16    return acc

```

## ✓ Training CoOpCLIP on Base Class

- The **Base class** refers to classes that were seen during training.
- In contrast, the **New class** refers to classes that were not seen during training.

```

1 # Train on the Base Classes Train split and evaluate accuracy on the Base Classes Test split.
2 args.trainer = "CoOp"
3 args.train_batch_size = 4
4 args.epoch = 100
5 args.output_dir = "outputs/coop"
6
7 args.subsample_classes = "base"
8 coop_base_acc = main(args)

```



```

epoch [94/100] batch [20/20] time 0.030 (0.060) data 0.000 (0.018) loss 0.0746 (0.0593) acc 100.0000 (100.0000) lr 3.0104e-05 eta 0:00:07
epoch [95/100] batch [20/20] time 0.031 (0.059) data 0.000 (0.016) loss 0.0079 (0.0987) acc 100.0000 (97.5000) lr 2.2141e-05 eta 0:00:05
epoch [96/100] batch [20/20] time 0.030 (0.060) data 0.000 (0.018) loss 0.0408 (0.0621) acc 100.0000 (100.0000) lr 1.5390e-05 eta 0:00:04
epoch [97/100] batch [20/20] time 0.032 (0.059) data 0.000 (0.016) loss 0.0678 (0.0630) acc 100.0000 (100.0000) lr 9.8566e-06 eta 0:00:03
epoch [98/100] batch [20/20] time 0.030 (0.059) data 0.000 (0.018) loss 0.0595 (0.0781) acc 100.0000 (98.7500) lr 5.5475e-06 eta 0:00:02
epoch [99/100] batch [20/20] time 0.032 (0.059) data 0.000 (0.019) loss 0.0272 (0.0819) acc 100.0000 (98.7500) lr 2.4666e-06 eta 0:00:01
epoch [100/100] batch [20/20] time 0.030 (0.059) data 0.000 (0.016) loss 0.0029 (0.0881) acc 100.0000 (97.5000) lr 6.1680e-07 eta 0:00:00
Checkpoint saved to outputs/coop/prompt_learner/model.pth.tar-100
Finish training
Deploy the last-epoch model
Evaluate on the *test* set
100%|██████████| 42/42 [00:19<00:00, 2.18it/s]=> result
* total: 4,200
* correct: 3,839
* accuracy: 91.4%
* error: 8.6%
* macro_f1: 91.5%
Elapsed: 0:02:53

```

## ▼ Evaluate CoOpCLIP on New Class

```

1 # Accuracy on the New Classes.
2 args.model_dir = "outputs/coop"
3 args.output_dir = "outputs/coop/new_classes"
4 args.subsample_classes = "new"
5 args.load_epoch = 100
6 args.eval_only = True
7 coop_novel_acc = main(args)

➡ Loading trainer: CoOp
Loading dataset: EuroSAT
Reading split from /content/ProMetaR/data/eurosat/split_zhou_EuroSAT.json
Loading preprocessed few-shot data from /content/ProMetaR/data/eurosat/split_fewshot/shot_16-seed_1.pkl
SUBSAMPLE NEW CLASSES!
Building transform_train
+ random resized crop (size=(224, 224), scale=(0.08, 1.0))
+ random flip
+ to torch tensor of range [0, 1]
+ normalization (mean=[0.48145466, 0.4578275, 0.40821073], std=[0.26862954, 0.26130258, 0.27577711])
Building transform_test
+ resize the smaller edge to 224
+ 224x224 center crop
+ to torch tensor of range [0, 1]
+ normalization (mean=[0.48145466, 0.4578275, 0.40821073], std=[0.26862954, 0.26130258, 0.27577711])

-----
Dataset      EuroSAT
# classes    5
# train_x    80
# val        20
# test       3,900
-----

Loading CLIP (backbone: ViT-B/16)
/usr/local/lib/python3.10/dist-packages/torch/utils/data/dataloader.py:617: UserWarning: This DataLoader will create 8 worker processes in total. Our
warnings.warn(
/usr/local/lib/python3.10/dist-packages/torch/optim/lr_scheduler.py:62: UserWarning: The verbose parameter is deprecated. Please use get_last_lr() to
warnings.warn(
/content/ProMetaR/dassl/utils/torchtools.py:102: FutureWarning: You are using `torch.load` with `weights_only=False` (the current default value), whic
checkpoint = torch.load(fpath, map_location=map_location)
Building custom CLIP
Initializing a generic context
Initial context: "X X X X X X X X X X X X X X X X"
Number of context words (tokens): 16
Turning off gradients in both the image and the text encoder
Loading evaluator: Classification
Loading weights to prompt_learner from "outputs/coop/prompt_learner/model.pth.tar-100" (epoch = 100)
Evaluate on the *test* set
100%|██████████| 39/39 [00:15<00:00, 2.59it/s]=> result
* total: 3,900
* correct: 2,007
* accuracy: 51.5%
* error: 48.5%
* macro_f1: 45.6%

```

## ▼ Visualization

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 metrics = ['Base', 'Novel']
5
6 coop_acc_list = [coop_base_acc, coop_novel_acc]
7
8 font_size = 12

```

```
8 bar_width = 0.33
9 index = np.arange(len(metrics))
10 fig, ax = plt.subplots()
11 bar1 = ax.bar(index, coop_acc_list, bar_width, label='CoOp')
12
13 ax.set_ylabel('Scores')
14 ax.set_title('Model Performance Comparison')
15 ax.set_xticks(index + bar_width / 2)
16 ax.set_xticklabels(metrics)
17 ax.legend()
18
19 def add_value_labels(bars):
20     for bar in bars:
21         height = bar.get_height()
22         ax.annotate(f'{height:.2f}', xy=(bar.get_x() + bar.get_width() / 2, height),
23                   xytext=(0, 2), # 2 points vertical offset
24                   textcoords='offset points',
25                   ha='center', va='bottom')
26
27 add_value_labels(bar1)
28 plt.tight_layout()
29 plt.show()
30
```

