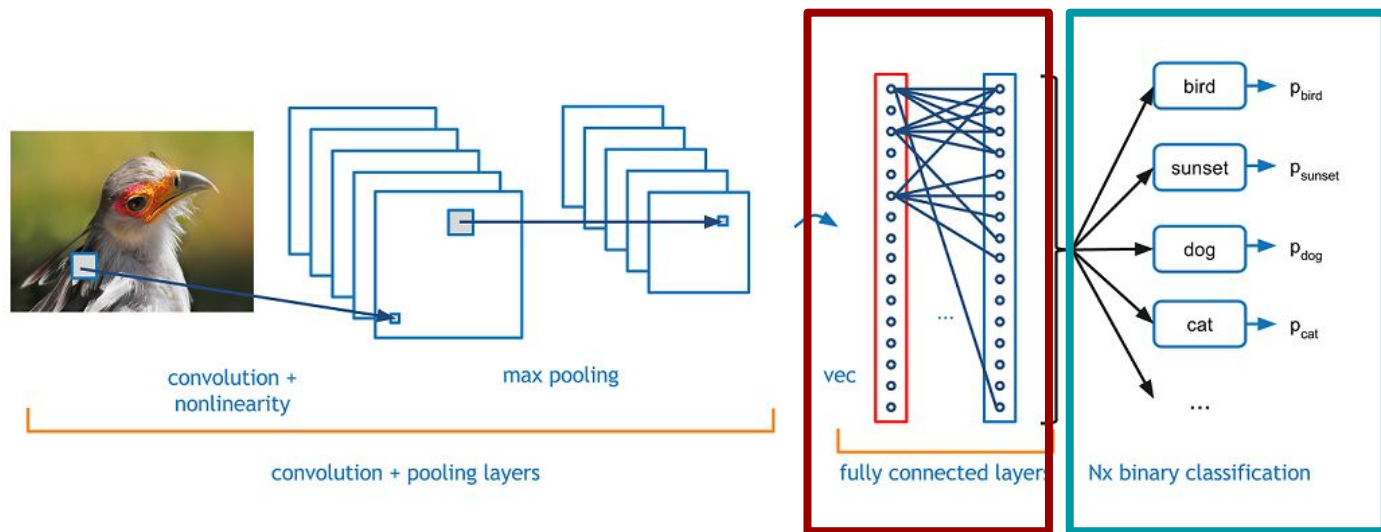


Land Cover Classification Using Foundation Models

Haiyang Jiang, Jingnan Cao

Introduction

Traditional Vision System Limitations



Predefined classifier

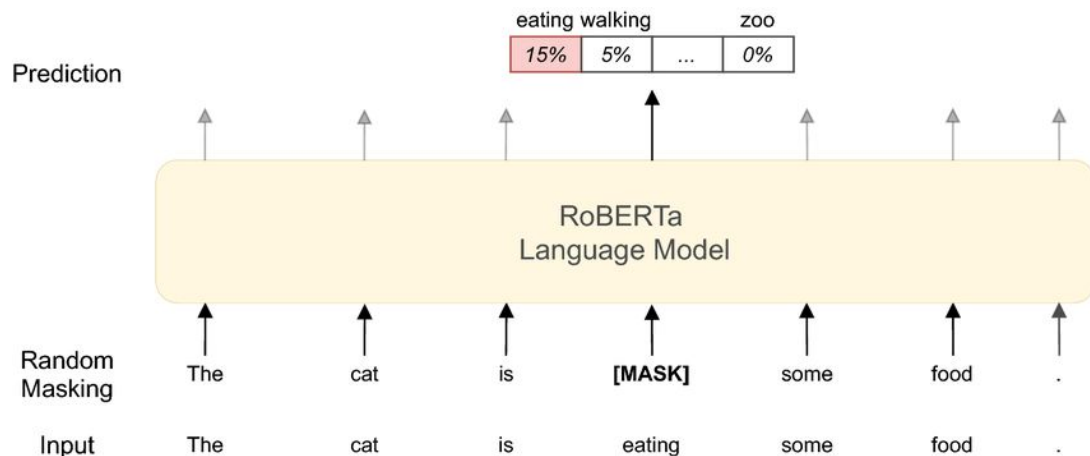
- > Static output
- > Lack of dynamic

Predefined object classes

- > Fixed set solution
- > Hard to add additional label
- > Limited generability

NLP Revolution by BERT & GPT

Directly learn from **raw text** by task-agnostic objectives
(autoregressive / masked language modeling)



- > Can learn from rich online resources -> **Scalability**
- > Boarder source of supervision from the text -> **Self-supervised**
- > Ability to learn **generalized** feature representation
- > Strong zero-shot **generalization** to downstream tasks

Language vs. Vision

Natural Language

- Self supervision (LM)
- Large training data
- Zero-shot transferability



Images

- Supervised learning
- Not that large training data (ImageNet)



Idea: enabling better transferability by connecting vision tasks by languages guiding

Related Work

Learning Image Representation from Natural Language

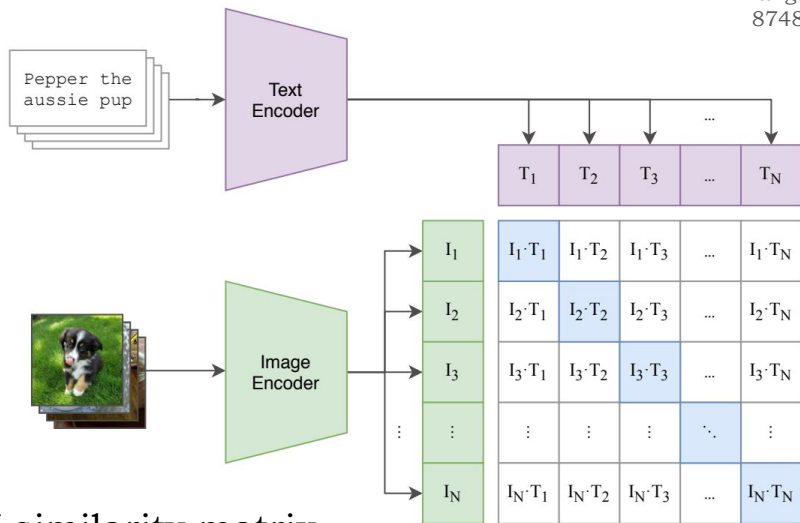
By reviewing historical works, reasons why this idea haven't achieved impressive performance:

1. Image-text paired dataset not large enough
 - > WIT WebImageText: **400M image-text pairs**
 - > Comparable with WebText for GPT-2
2. Model not powerful enough
 - > **ResNet / Vision Transformer** as vision encoder
 - > **Transformer** as text encoder

CLIP: Contrastive Language-Image Pretraining

Train on WebImageText(WIT): a newly constructed dataset of **400 million** (image, text) pairs on the Internet

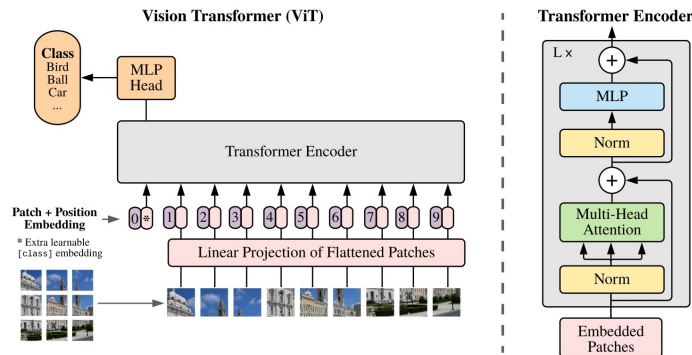
Radford A, Kim J W, Hallacy C, et al. Learning transferable visual models from natural language supervision[C]//International conference on machine learning. PMLR, 2021: 8748-8763.



$N \times N$ similarity matrix

Maximizing **N** pairs similarity (positive pair)

Minimizing $N^2 - N$ pairs similarity (negative pair)



$$\min \left(\sum_{i=1}^N \sum_{j=1}^N (I_i \cdot T_j)_{i \neq j} - \sum_{i=1}^N (I_i \cdot T_i) \right)$$

batch-size = 32,768 **VERY BIG!**

Training Objective

Crucial Problem: for scaled data,

training efficiency is the key to success

Start from predicting **caption** of the image (similar in generative model)

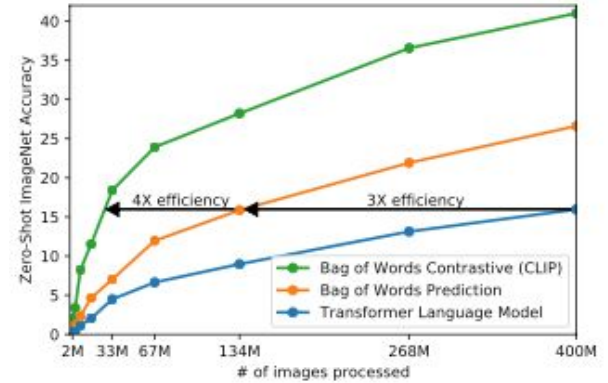
-> Wide variety of possible answer -> **Hard & high computation cost**

Captioning: question answering

Classification: multiple choice

Pairing: simple yes or no

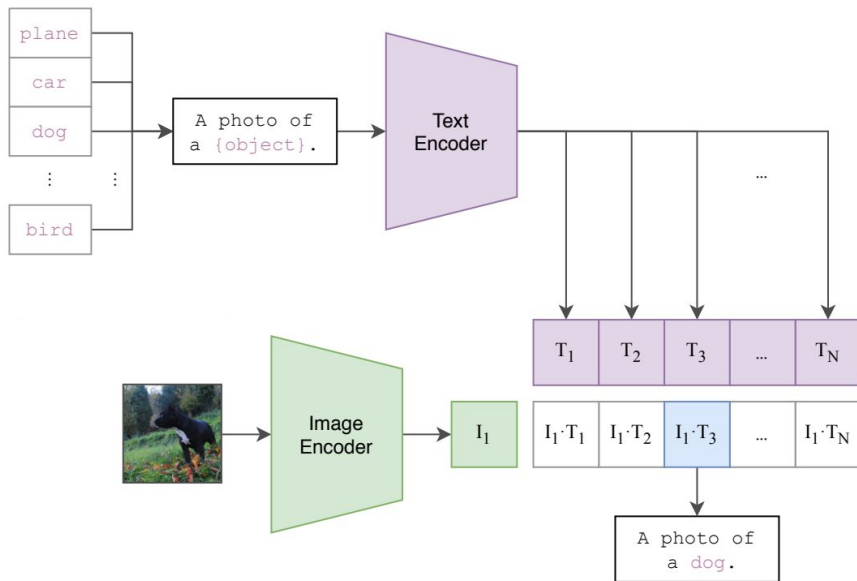
-> **Contrastive learning**



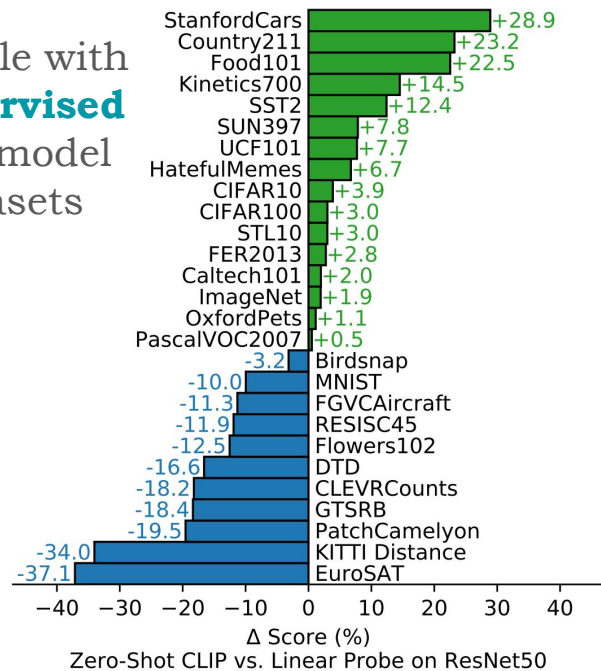
**computation
cost reducing**



Zero-Shot Image Classification



Comparable with
fully supervised
ResNet50 model
on 27 datasets



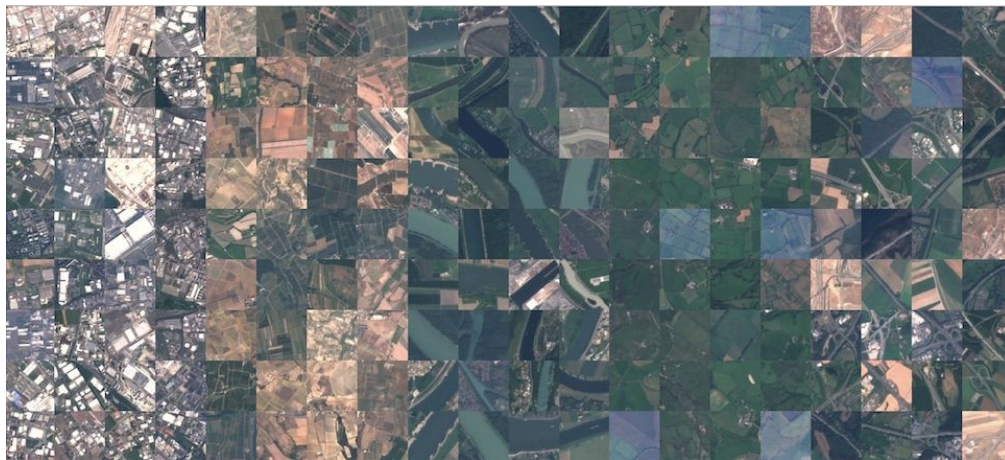
Class + Template = Sentence

Each image \rightarrow N classes similarity \rightarrow Get the max as prediction

Task

Zero-shot / Few-shot Classification on EuroSAT

Using Foundation Model CLIP for land cover classification on the EuroSAT dataset to classify satellite images into various land cover types



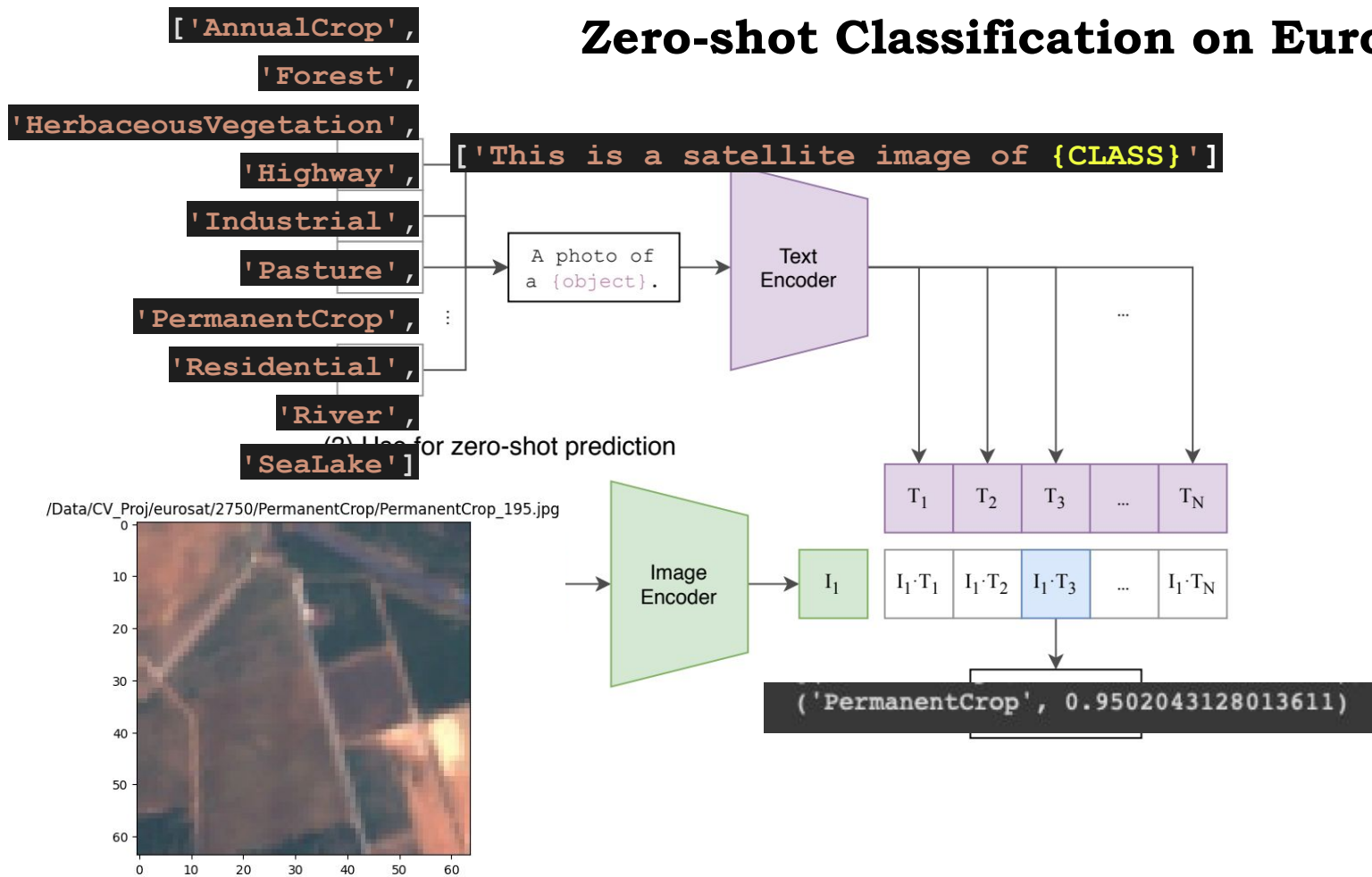
EuroSAT covers **13 spectral bands** and consisting out of **10 classes** with in total **27,000 labeled** and geo-referenced images

```
['AnnualCrop', 'Forest', 'HerbaceousVegetation', 'Highway',
```

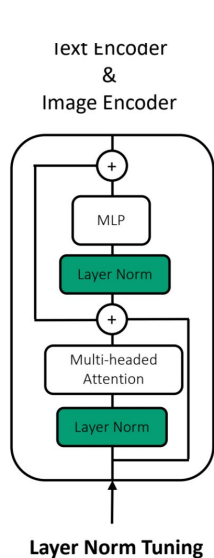
```
'Industrial', 'Pasture', 'PermanentCrop', 'Residential', 'River', 'SeaLake']
```

Method

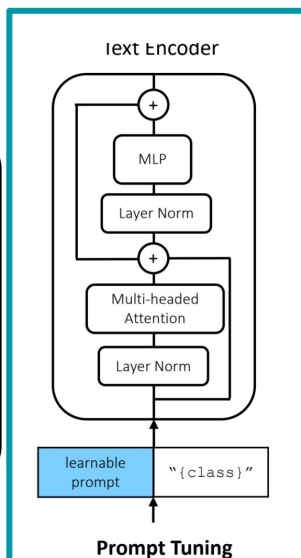
Zero-shot Classification on EuroSAT



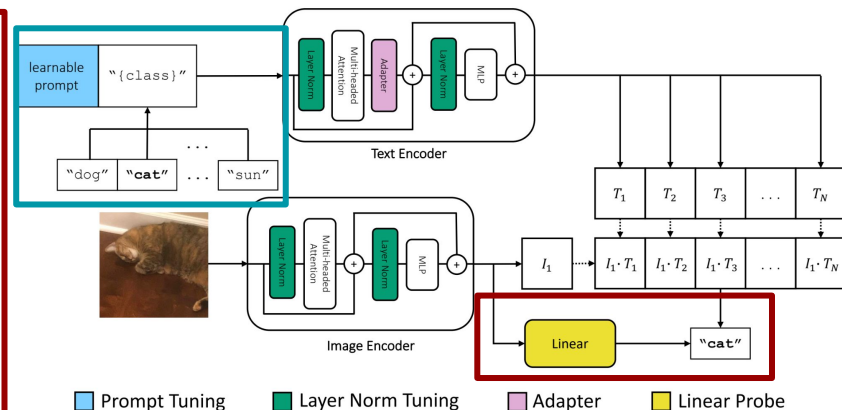
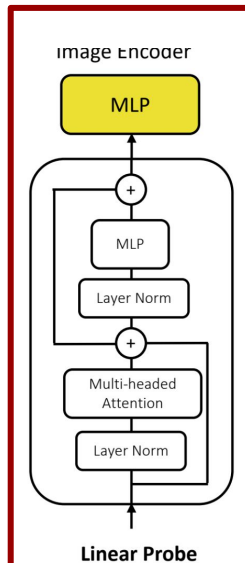
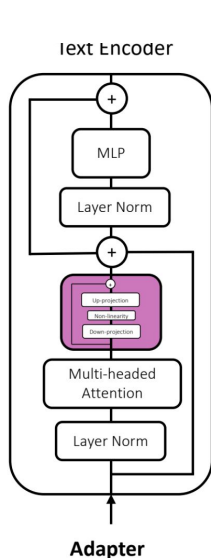
Few-Shot Learning Enhancement



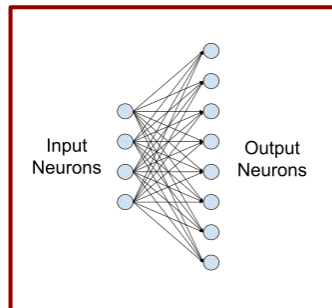
No possible for Foundation Model !



No possible for Foundation Model !

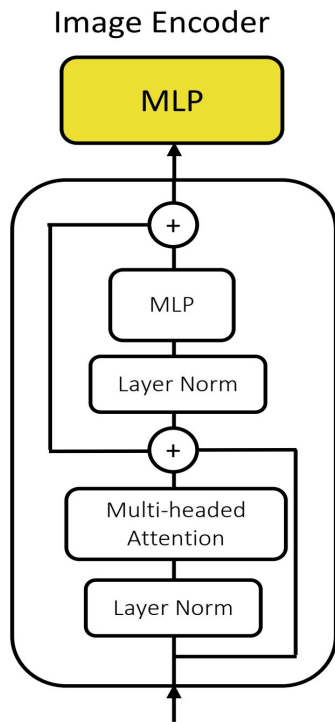


```
imagenet_templates = [
    'a bad photo of a {}.',
    'a photo of many {}.',
    'a sculpture of a {}.',
    'a photo of the hard to see {}.',
    'a low resolution photo of the {}.',
    'a rendering of a {}.',
    'graffiti of a {}.',
    'a bad photo of the {}.',
    'a cropped photo of the {}.',
    'a tattoo of a {}.',
    'the embroidered {}.',
    'a photo of a hard to see {}.',
    'a bright photo of a {}.',
    'a photo of a clean {}.',
    'a photo of a dirty {}.',
    'a dark photo of the {}.',
    'a drawing of a {}.',
    'a photo of my {}.',
]
```



Few-Shot Learning Enhancement

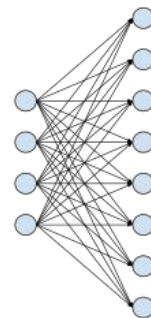
- Linear Probe



1. 16 shots learning
16 samples / class for training
2. Full dataset feature learning
7:3 Training vs. Validating

Both test on whole dataset

dimension
of vision
encoding



number of
classes

```
classifier = SimpleClassifier(input_dim, num_classes).to(device)
optimizer = optim.Adam(classifier.parameters(), lr=1e-3)
loss_fn = nn.CrossEntropyLoss()

num_epochs = 100
validation_interval = 10 # Validate after every 10 epochs
save_interval = 20 # Save after every 20 epochs
early_stopping_patience = 10 # Number of epochs to wait for loss improvement
best_val_accuracy = 0 # Track best validation accuracy
no_improvement_epochs = 0 # Track epochs without improvement for early stopping
```

Few-Shot Learning Enhancement

- Prompt Learning CoOp

Hand-crafted prompt engineering limitation:

1. Require domain expertise
2. Extremely time-consuming
3. Slight change will cause big difference

Automatically learn the contextual vectors in continuous space

-> Use cross-entropy classification loss to evaluate a context vector towards a class

-> Gradient from text encoder all the way to original context vectors

-> Data-efficient, beat hand-crafted prompts & linear probe with 1-2 shots

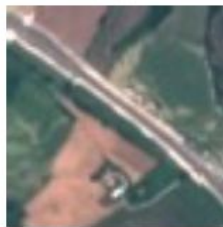
Flowers102



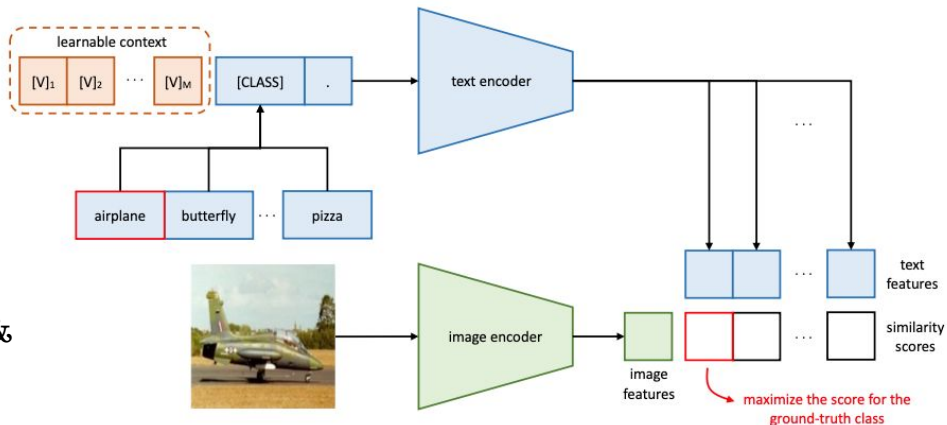
Prompt	Accuracy
a photo of a [CLASS].	60.86
a flower photo of a [CLASS].	65.81
a photo of a [CLASS], a type of flower.	66.14
$[V]_1 [V]_2 \dots [V]_M$ [CLASS].	94.51

(b)

EuroSAT



Prompt	Accuracy
a photo of a [CLASS].	24.17
a satellite photo of [CLASS].	37.46
a centered satellite photo of [CLASS].	37.56
$[V]_1 [V]_2 \dots [V]_M$ [CLASS].	83.53



Few-Shot Learning Enhancement

- Prompt Learning CoOp & CoCoOp

$$p(y = i | \mathbf{x}) = \frac{\exp(\cos(g(\mathbf{t}_i), \mathbf{f})/\tau)}{\sum_{j=1}^K \exp(\cos(g(\mathbf{t}_j), \mathbf{f})/\tau)},$$

Simple cross entropy objective

$$\mathbf{t} = [\mathbf{V}]_1 [\mathbf{V}]_2 \dots [\mathbf{V}]_M [\text{CLASS}],$$

$$\mathbf{t}_i = [\mathbf{V}]_1^i [\mathbf{V}]_2^i \dots [\mathbf{V}]_M^i [\text{CLASS}]^i,$$

$$\mathbf{t}_i(\mathbf{x}) = \{v_1(\mathbf{x}), v_2(\mathbf{x}), \dots, v_M(\mathbf{x}), c_i\}$$

$$\pi = h_{\theta}(\mathbf{x}) \quad \text{image-based token}$$

$$v_m(\mathbf{x}) = v_m + \pi \quad \text{merge image token and context vector}$$

1. Unified context for all classes

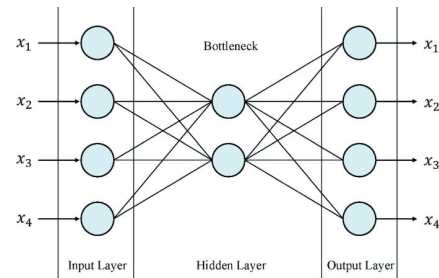
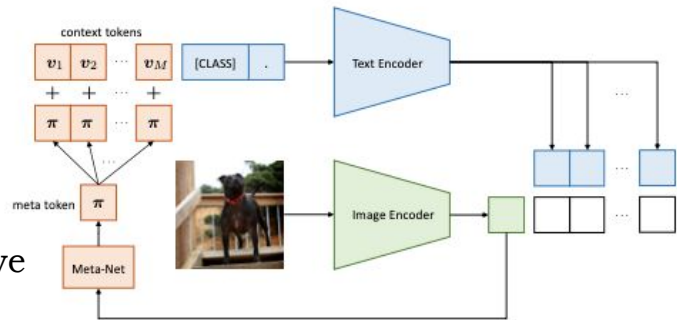
-> More generalized context
2. Class-specified context

-> might benefit fine-grained dataset such as StanfordCars
 -> more parameters -> need more shots
 -> limited generalization to new class

3. Context conditioned on image

Above prompt: poor performance on unseen class -> static context

-> from class-oriented to instance-oriented
 -> dynamic context based on image
 -> better generalized on unseen class
 -> much more parameters to train



Few-Shot Learning Enhancement

- Prompt Learning CoOp & CoCoOp

$$t = [V]_1[V]_2 \dots [V]_M[\text{CLASS}],$$

1. Unified context for all classes

16 context + 1 class = 17 words

-> learn tensor [1*16*77]

$$t_i = [V]_1^i[V]_2^i \dots [V]_M^i[\text{CLASS}]^i,$$

2. Class-specified context

(16 context + 1 class) * 10 classes = 170 words

-> learn tensor [10*16*77]

$$t_i(x) = \{v_1(x), v_2(x), \dots, v_M(x), c_i\}$$

$$\pi = h_{\theta}(x)$$

image-based token

$$v_m(x) = v_m + \pi$$

merge image token
and context vector

3. Context conditioned on image

Starts from initialized “a photo of a {CLASS}”

4 context + 1 class = 5 words

-> learn tensor [1*4*77]

+

Meta-Net (two-layer bottleneck)

**Since introduce a network, to
keep training efficiency reduce
the context length to 4**

Result

Zero-shot CLIP

1849	48	1	261	233	49	487	43	3	26
17	1072	0	1206	16	286	5	6	8	384
51	644	2	601	20	841	25	66	16	734
323	8	0	1598	235	116	138	61	5	16
48	0	0	20	1813	6	1	612	0	0
499	32	0	25	26	1142	234	4	0	0
1071	8	1	156	407	132	520	170	1	34
18	5	0	43	280	5	0	2647	1	1
328	12	0	473	298	258	119	50	371	591
4	916	0	184	44	13	3	3	1228	605

F1-Score:

43%

-> Not good

16-shots Linear Probe

2460	10	28	27	0	53	385	0	6	31
1	2363	93	3	0	76	0	1	2	461
2	159	2615	78	2	30	48	10	11	15
225	31	31	1739	83	117	127	54	93	6
35	1	5	13	2334	0	25	86	1	0
111	100	61	26	2	1582	105	0	7	6
363	3	156	113	65	140	1583	69	4	4
5	14	32	4	160	0	4	2779	1	1
177	57	41	279	25	99	58	7	1744	13
8	131	30	5	0	2	1	1	88	2734

F1-Score:

81% (+38%)

-> Huge

improve

-> Close to upper limit

Representation Learning Full Dataset Linear Probe

563	0	1	4	1	6	12	0	6	1
0	576	6	0	0	2	0	0	1	2
1	12	588	1	0	2	10	0	2	1
10	5	3	431	4	2	8	3	11	0
0	0	0	6	513	0	1	5	1	0
8	7	10	2	0	374	3	0	2	0
27	0	15	5	2	3	419	0	2	0
0	1	0	1	5	0	4	595	0	0
13	1	2	27	4	5	0	0	461	3
1	2	2	1	0	0	0	0	9	581

F1-Score:

94% (+51%)

-> Upper limit

16-shots Unified Context Prompt Learning + Zero-shot CLIP

2295	5	7	121	0	104	251	0	204	13
1	2834	41	5	0	61	10	0	3	45
0	317	2227	37	2	84	297	5	7	24
58	12	31	1722	112	98	147	36	284	0
2	0	0	6	2417	1	52	18	4	0
1	37	23	54	0	1716	132	1	36	0
70	2	52	70	59	26	2201	5	15	10
1	18	12	5	121	5	88	2742	4	4
27	37	27	377	58	176	47	8	1737	6
3	43	11	3	0	87	1	0	133	2719

F1-Score:

83% (+40%)

-> Comparable linear probe

-> Close to upper limit

16-shots Class-specified Context Prompt Learning + Zero-shot CLIP

2232	2	12	178	0	101	193	0	269	13
1	2894	31	3	0	25	2	1	5	38
1	352	2213	29	2	67	295	11	10	20
52	9	48	1680	92	127	136	59	296	1
2	0	0	11	2395	2	43	43	4	0
1	57	11	16	0	1818	78	1	16	0
57	2	44	52	49	26	2223	16	31	0
2	15	10	9	92	2	66	2803	1	0
17	16	48	421	49	178	30	20	1711	10
6	14	2	3	1	72	0	0	128	2774

F1-Score:

84% (+41%)

-> Comparable linear probe

-> Close to upper limit

-> No obvious advantage on class-specified

16-shots Conditional Context Prompt Learning + Zero-shot CLIP

Train on 5 classes (1/2)
Test on all 10 classes

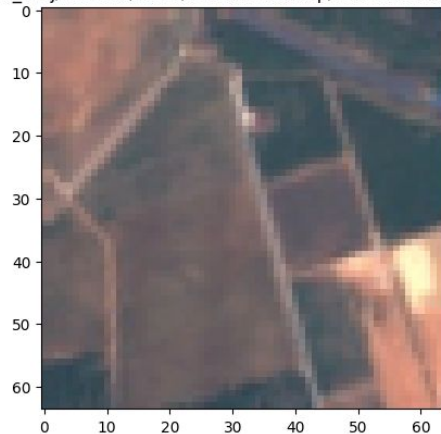
F1-Score:

65% (+22%)

-> Improved generalization compared to baseline

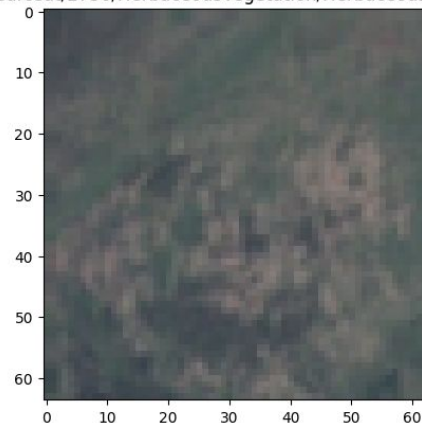
Some Results from Unified Context Prompt

/Data/CV_Proj/eurosat/2750/PermanentCrop/PermanentCrop_195.jpg



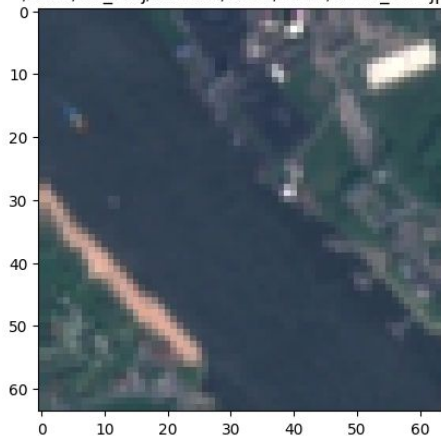
```
[('AnnualCrop', 0.0112),  
 ('Forest', 4e-05),  
 ('HerbaceousVegetation',  
 0.0004),  
 ('Highway', 0.0322),  
 ('Industrial', 0.0018),  
 ('Pasture', 0.0036),  
 ('PermanentCrop', 0.9502),  
 ('Residential', 0.0002),  
 ('River', 0.0004),  
 ('SeaLake', 5e-06)]
```

/Data/CV_Proj/eurosat/2750/HerbaceousVegetation/HerbaceousVegetation_2069.jpg



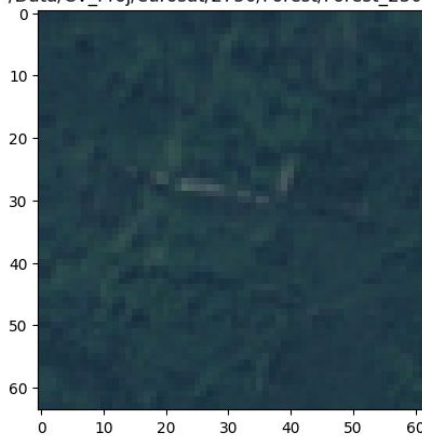
```
[('AnnualCrop', 1e-05),  
 ('Forest', 0.1713),  
 ('HerbaceousVegetation',  
 0.7921),  
 ('Highway', 0.0010),  
 ('Industrial', 0.0002),  
 ('Pasture', 0.0147),  
 ('PermanentCrop', 0.0164),  
 ('Residential', 0.0030),  
 ('River', 0.0001),  
 ('SeaLake', 0.0012)]
```

/Data/CV_Proj/eurosat/2750/River/River_201.jpg



```
[('AnnualCrop', 0.0002),  
 ('Forest', 0.0032),  
 ('HerbaceousVegetation',  
 0.0042),  
 ('Highway', 0.0185),  
 ('Industrial', 0.0196),  
 ('Pasture', 0.0365),  
 ('PermanentCrop', 0.0053),  
 ('Residential', 0.0078),  
 ('River', 0.9032),  
 ('SeaLake', 0.0014)]
```

/Data/CV_Proj/eurosat/2750/Forest/Forest_2303.jpg



```
[('AnnualCrop', 0.0001),  
 ('Forest', 0.7211),  
 ('HerbaceousVegetation',  
 0.0647),  
 ('Highway', 0.0018),  
 ('Industrial', 0.0424),  
 ('Pasture', 0.0999),  
 ('PermanentCrop', 0.0510),  
 ('Residential', 0.0171),  
 ('River', 0.0008),  
 ('SeaLake', 0.0012)]
```

Interpret the Learned Prompt

from Unified Context Prompt

- 1: [**'decorations'**</w>, 'lizards</w>', 'wed</w>', 'dor', 'erin</w>'] [0.6720, '0.6734', '0.6743', '0.6754', '0.6755']
- 2: [**'pelo'**, 'sculpted</w>', 'lit</w>', 'revol', 'appeared</w>'] [0.8290, '0.8376', '0.8379', '0.8389', '0.8390']
- 3: ['jake', 'jw', 'joe', 'kab</w>', 'half'] [0.7303, '0.7352', '0.7359', '0.7373', '0.7385']
- 4: ['blames</w>', **'organised'**</w>', 'applaud</w>', 'picked</w>', 'implic'] [0.8232, '0.8240', '0.8244', '0.8251', '0.8258']
- 5: ['list', **'rotating'**</w>', **'represented'**</w>', 'strack</w>', 'alline</w>'] [0.7324, '0.7375', '0.7377', '0.7402', '0.7413']
- 6: ['knife</w>', 'etu', 'broken', **'foreign'**, 'exploding</w>'] [1.0701, '1.0703', '1.0724', '1.0741', '1.0744']
- 7: ['sweat', 'stri', 'tall', 'masa', 'yaw'] [0.9688, '0.9709', '0.9724', '0.9756', '0.9761']
- 8: [**'linear'**</w>', 'optional</w>', 'logical</w>', 'smear</w>', 'phillips</w>'] [0.8934, '0.8982', '0.9033', '0.9046', '0.9047']
- 9: ['piss</w>', 'simplified</w>', 'kow</w>', **'modes'**</w>', **'calm'**] [0.8488, '0.8512', '0.8536', '0.8564', '0.8569']
- 10: ['terri', 'ting', 'newsp', 'cops</w>', **'relocated'**</w>] [0.9932, '0.9934', '0.9936', '0.9950', '0.9950']
- 11: ['milb</w>', 'taxpayer</w>', ':|</w>', 'moms', 'nsfw</w>'] [0.9133, '0.9133', '0.9152', '0.9166', '0.9181']
- 12: ['tummy</w>', **'residence'**</w>', **'retreat'**</w>', 'chest</w>', 'anger</w>'] [0.8954, '0.8981', '0.8997', '0.9000', '0.9003']
- 13: ['tivity</w>', 'dar</w>', 'ities</w>', **'wood'**</w>', **'sight'**</w>] [0.6240, '0.6274', '0.6319', '0.6323', '0.6347']
- 14: [**'rooms'**</w>', 'officer</w>', **'weather'**</w>', **'toward'**</w>', **'lighting'**</w>] [0.6663, '0.6685', '0.6696', '0.6725', '0.6729']
- 15: ['ilove', 'ignored</w>', 'rig</w>', 'vi</w>', 'gamer</w>'] [0.7930, '0.7935', '0.7953', '0.7959', '0.7962']
- 16: [**'addresses'**</w>', **'cyber'**</w>', 'aring</w>', '', 'irl</w>'] [0.8209, '0.8216', '0.8233', '0.8234', '0.8237']

Problem with continuous
prompt vector learning:
Hard to interpret the result

Searching within the
vocabulary for words that
closest to the learned vector,
by Euclidean distance

Conclusion

Conclusion

1. CLIP zero-shot on EuroSAT is not good (actually in the paper, among the 27 datasets, EuroSAT is one of the hardest to CLIP)
2. Linear probe can improve the performance significantly with only 16 shots, which proves CLIP already learn the generalized image representation
3. Fully train the linear head on the whole dataset can push the performance to nearly 95%
-> Fully capable for potential downstream task
4. Prompt learning is an auxiliary approach to deploy CLIP to downstream task, performance similar to linear probe using same shots learning

Advantage to linear probe:

- a. More dynamic output and generalization potential
- b. Better transparency and explainability

Reference

- [1] Radford A, Kim J W, Hallacy C, et al. Learning transferable visual models from natural language supervision[C]//International conference on machine learning. PMLR, 2021: 8748-8763.
- [2] Kim K, Laskin M, Mordatch I, et al. How to adapt your large-scale vision-and-language model[J]. 2021.
- [3] Zhou, Kaiyang, et al. "Learning to prompt for vision-language models." *International Journal of Computer Vision* 130.9 (2022): 2337-2348.
- [4] Zhou, Kaiyang, et al. "Conditional prompt learning for vision-language models." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2022.