# Semantic Prompt for Few-Shot Image Recognition

Wentao Chen[1,2]* Chenyang Si[3]* Zhang Zhang[2,4] Liang Wang[2,4] Zilei Wang[1] Tieniu Tan[1,2,4]

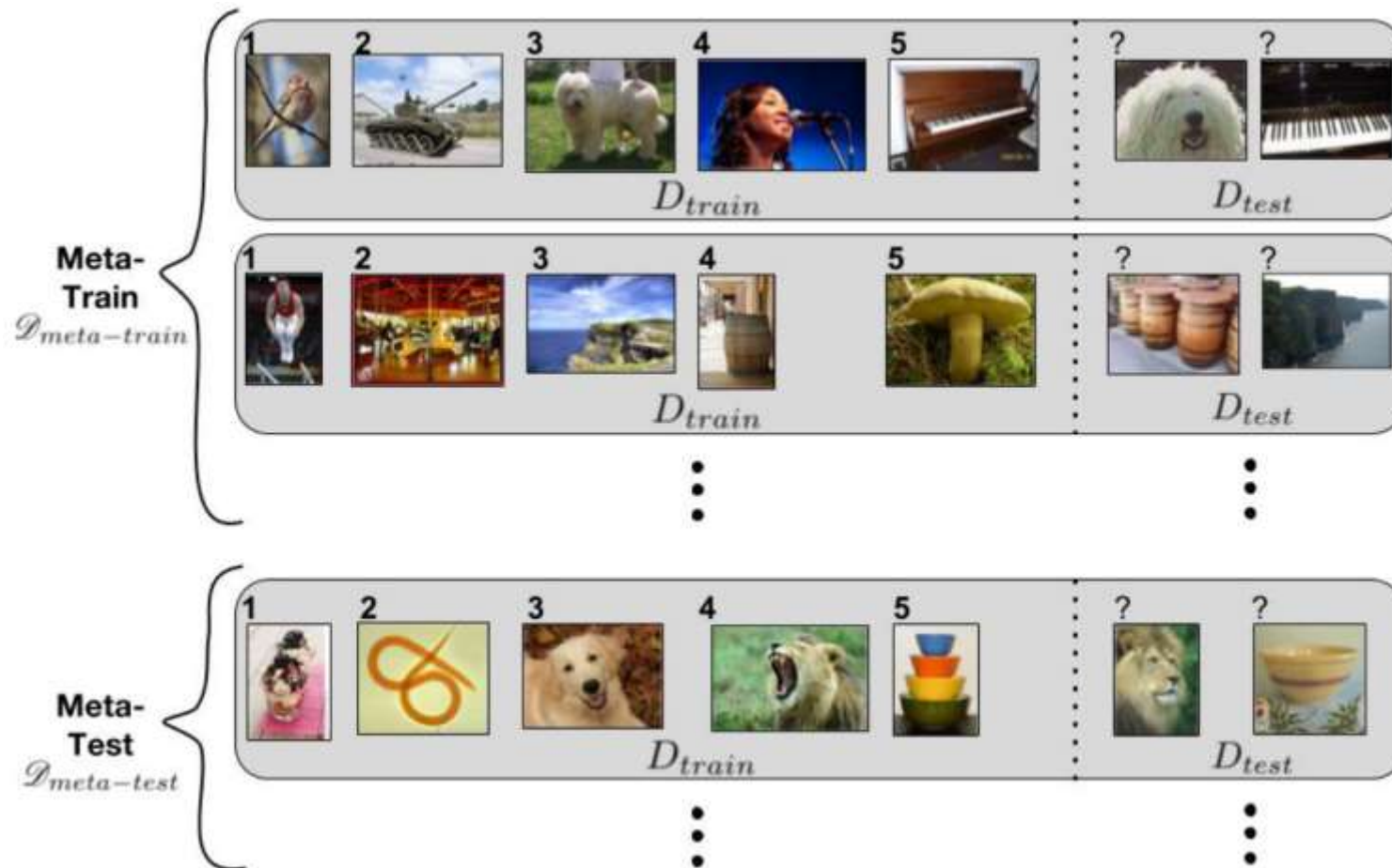[1]University of Science and Technology of China
[2]Center for Research on Intelligent Perception and Computing, NLPR, CASIA
[3]Nanyang Technological University [4]University of Chinese Academy of Sciences

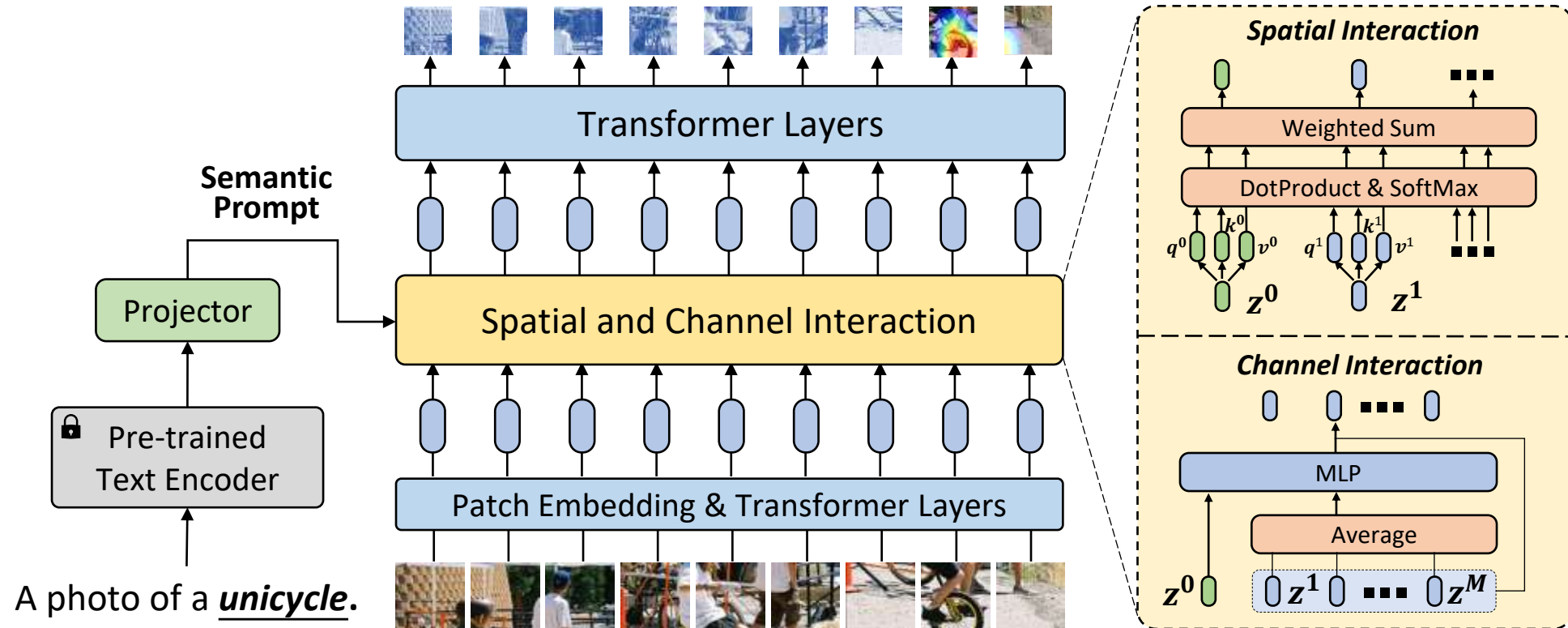THU-PM-284

# Quick Preview

■ Task

- We focus on the few-shot image recognition task, where only one or a few support images are available for a new class, and a large base dataset is used for meta-training.

# Quick Preview

■ Method

- We propose to use text data as semantic prompts to improve the visual feature extraction.
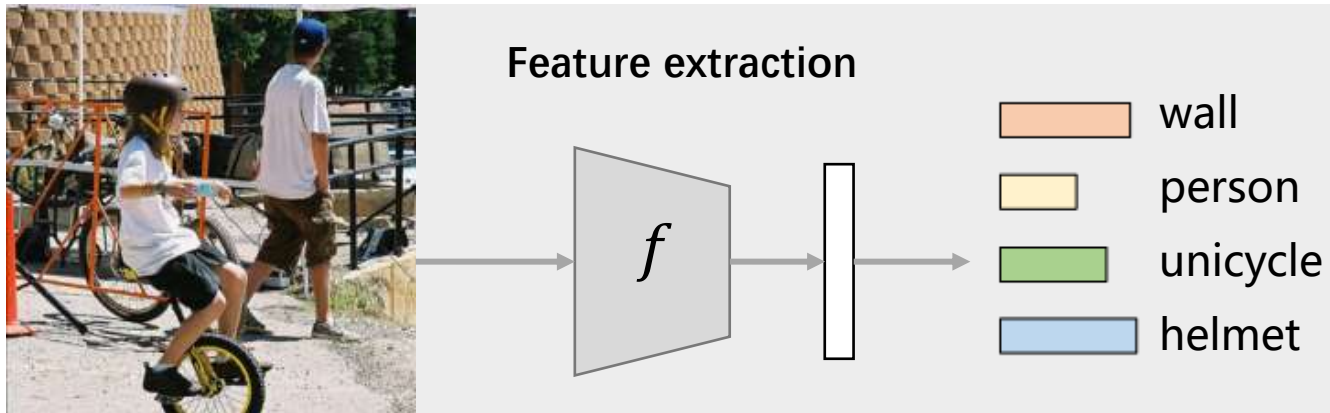
# Quick Preview

■ Experiments
- We evluate three different text encoders, and achieve consistent improvements on four datasets.

| Method | Backbone | Params/FLOPS | *mini*ImageNet 5-way | | *tiered*ImageNet 5-way | |
|--------|----------|--------------|---------|--------|---------|--------|
| | | | 1-shot | 5-shot | 1-shot | 5-shot |
| LEO [42] | WRN-28-10 | 36.5M/3.7 × 10¹⁰ | 61.76±0.08 | 77.59±0.12 | 66.33±0.05 | 81.44±0.09 |
| CC+rot [14] | WRN-28-10 | 36.5M/3.7 × 10¹⁰ | 62.93±0.45 | 79.87±0.33 | 70.53±0.51 | 84.98±0.36 |
| Align [1] | WRN-28-10 | 36.5M/3.7 × 10¹⁰ | 65.92±0.60 | 82.85±0.55 | **74.40±0.68** | 86.61±0.59 |
| MetaOptNet [22] | ResNet-12 | 12.5M/3.5 × 10⁹ | 62.64±0.61 | 78.63±0.46 | 65.99±0.72 | 81.56±0.53 |
| Meta-Baseline [6] | ResNet-12 | 12.5M/3.5 × 10⁹ | 63.17±0.23 | 79.26±0.17 | 68.62±0.27 | 83.74±0.18 |
| DeepEMD [56] | ResNet-12 | 12.5M/3.5 × 10⁹ | 65.91±0.82 | 82.41±0.56 | 71.16±0.87 | 86.03±0.58 |
| RE-Net [17] | ResNet-12 | 12.5M/3.5 × 10⁹ | 67.60±0.44 | 82.58±0.30 | 71.61±0.51 | 85.28±0.35 |
| TPMM [51] | ResNet-12 | 12.5M/3.5 × 10⁹ | 67.64±0.63 | **83.44±0.43** | 72.24±0.70 | 86.55±0.63 |
| SetFeat [2] | ResNet-12 | 12.5M/3.5 × 10⁹ | **68.32±0.62** | 82.71±0.46 | 73.63±0.88 | **87.59±0.57** |
| SUN [10] | Visformer-S | 12.4M/1.7 × 10⁸ | 67.80±0.45 | 83.25±0.30 | 72.99±0.50 | 86.74±0.33 |
| KTN [32] | ResNet-12 | 12.5M/3.5 × 10⁹ | 61.42±0.72 | 74.16±0.56 | - | - |
| AM3 [52] | ResNet-12 | 12.5M/3.5 × 10⁹ | 65.30±0.49 | 78.10±0.36 | 69.08±0.47 | 82.58±0.31 |
| TRAML [24] | ResNet-12 | 12.5M/3.5 × 10⁹ | **67.10±0.52** | 79.54±0.60 | - | - |
| DeepEMD-BERT [53] | ResNet-12 | 12.5M/3.5 × 10⁹ | 67.03±0.79 | **83.68±0.65** | **73.76±0.72** | **87.51±0.75** |
| Pre-train (Ours) | Visformer-T | 10.0M/1.3 × 10⁹ | 65.16±0.44 | 81.22±0.32 | 72.38±0.50 | 86.74±0.34 |
| **SP-CLIP** (Ours) | Visformer-T | 10.0M/1.3 × 10⁹ | **72.31±0.40** | 83.42±0.30 | **78.03±0.46** | 88.55±0.32 |
| **SP-SBERT** (Ours) | Visformer-T | 10.0M/1.3 × 10⁹ | 70.70±0.42 | **83.55±0.30** | 73.31±0.50 | 88.56±0.32 |
| **SP-GloVe** (Ours) | Visformer-T | 10.0M/1.3 × 10⁹ | 70.81±0.42 | 83.31±0.30 | 74.68±0.50 | **88.64±0.31** |

Table 1. Comparison with previous work on *mini*ImageNet and *tiered*ImageNet. Methods in the top rows do not use semantic information, and methods in the middle rows leverage semantic information from class names [24, 32, 52] or descriptions [53]. Accuracies are reported with 95% confidence intervals.

# Motivation

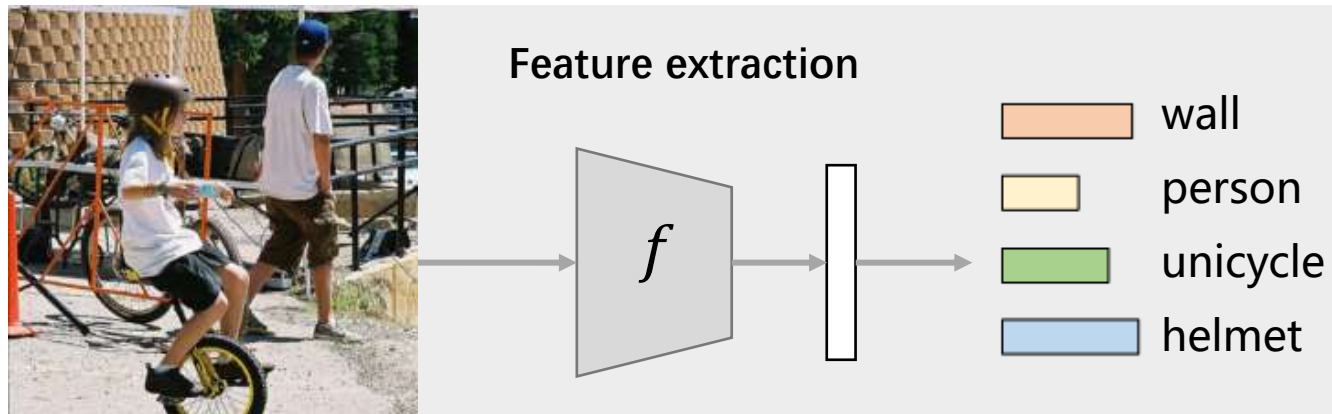- Given only one support image, the obtained image feature may contain much nioses.



Input image

{'unicycle'}

# Motivation

- Given only one support image, the obtained image feature may contain much nioses.
- The class name has rich semantic information that can be extracted by a text encoder.

# Motivation

- Given only one support image, the obtained image feature may contain much nioses.
- The class name has rich semantic information that can be extracted by a text encoder.
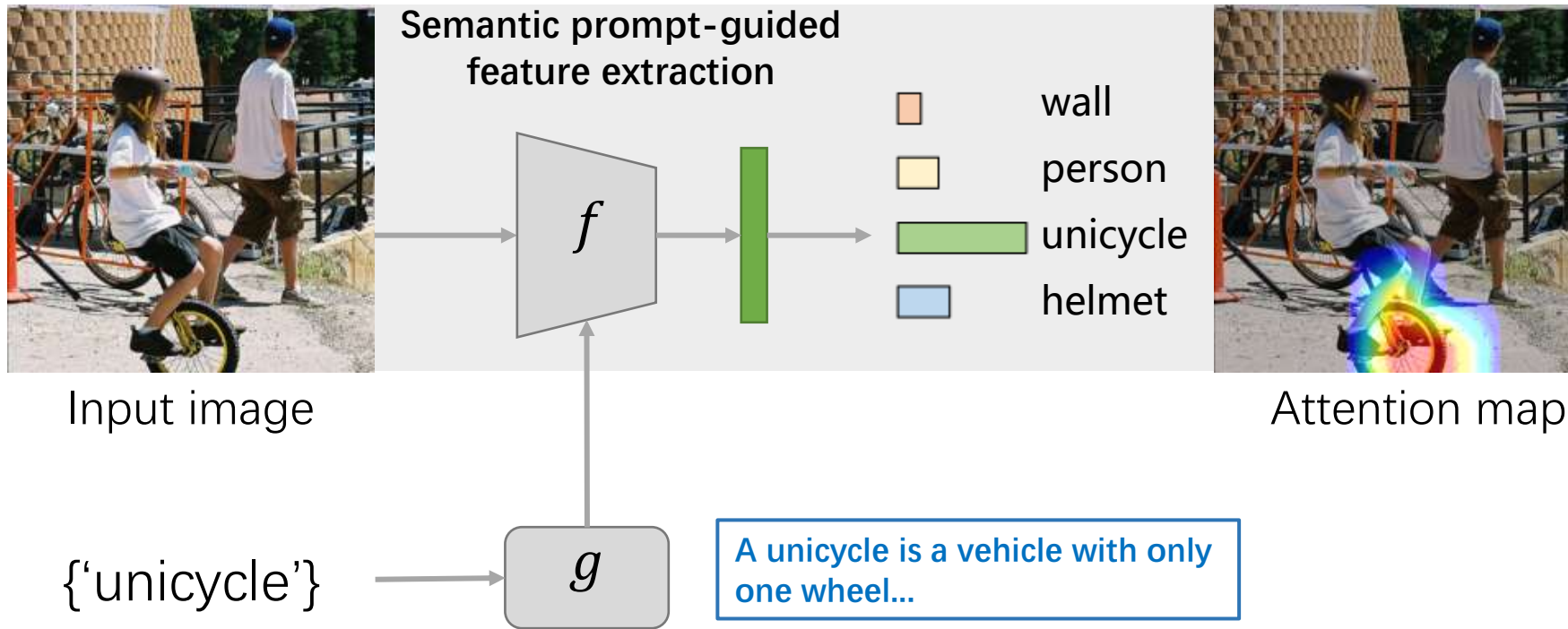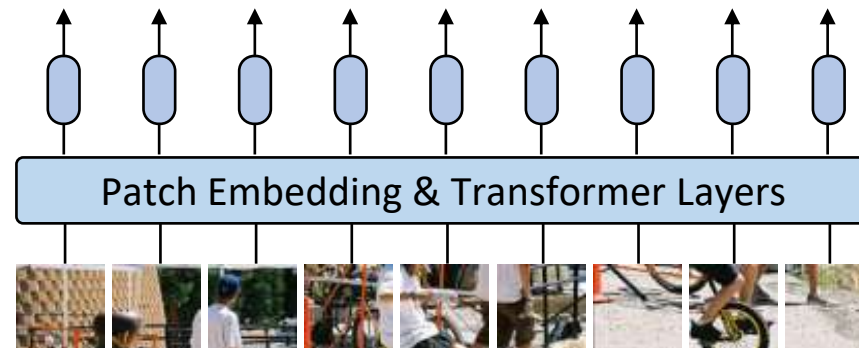- We use semantic features as prompts to improve the visual feature extraction.

# The framework of semantic prompt

- Feed image patches into a Vision Transformer.

# The framework of semantic prompt

- Feed image patches into a Vision Transformer.
- Feed the class name into a text encoder to obtain a semantic prompt.

# The framework of semantic prompt

- Feed image patches into a Vision Transformer.
- Feed the class name into a text encoder to obtain a semantic prompt.
- Extract image features guided by the semantic prompt via spatial and channel interaction.

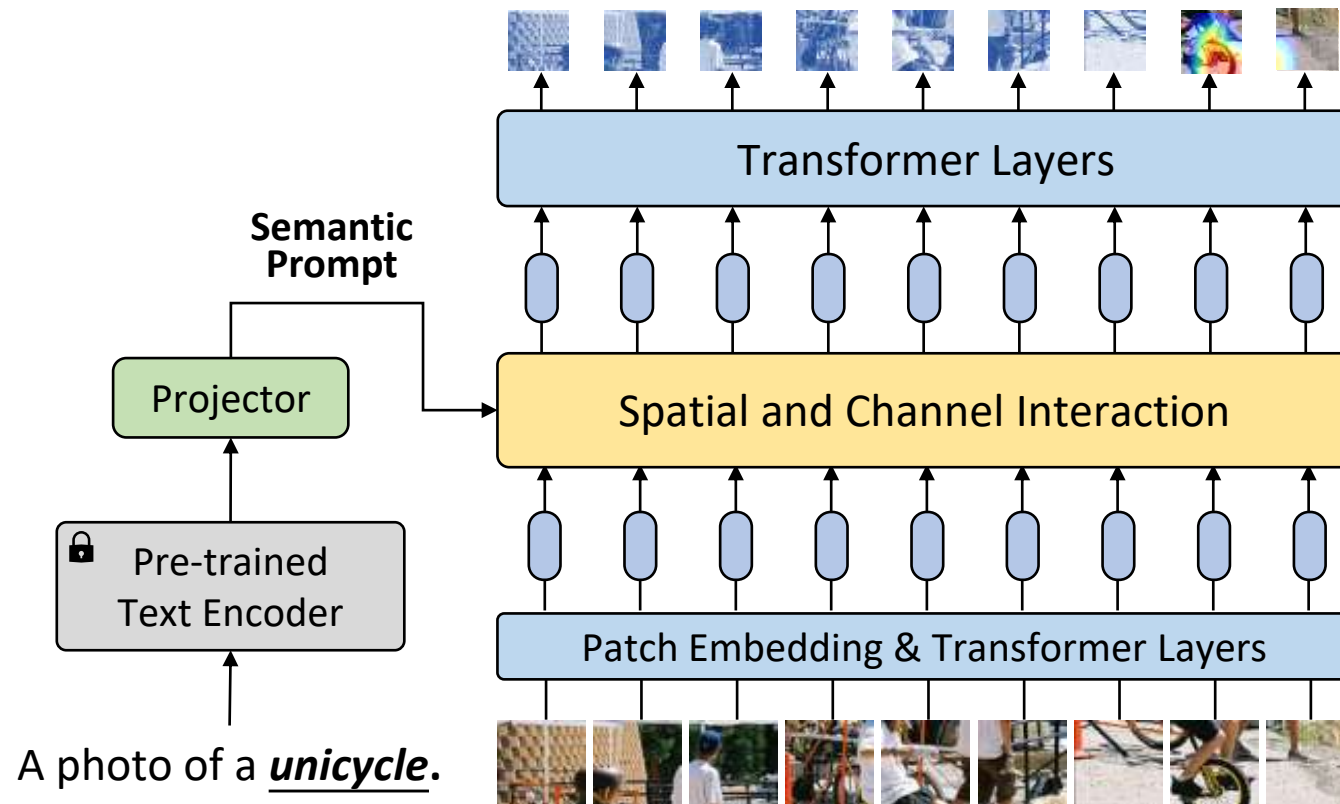# The framework of semantic prompt

- Feed image patches into a Vision Transformer.
- Feed the class name into a text encoder to obtain a semantic prompt.
- Extract image features guided by the semantic prompt via spatial and channel interaction.
- Train the model via meta-learning.



1.class prototype:

$$c_i = \frac{1}{K}\sum_{j=1}^{K} f_g(x_j^S)$$

2.loss function:

$$L_{meta} = -\mathbb{E}_{S,Q}\,\mathbb{E}_{x^q} \log \frac{\exp(s(f(x^q),c_{y^q})/\tau)}{\sum_{i=1}^{N}\exp(s(f(x^q),c_i)/\tau)}$$

# Spatial and channel interaction

- Adapt visual features on spatial and channel dimensions according to the given prompt.



- Spatial Interaction

  ① Concat the prompt and patches.
  $$\widehat{Z}_{l-1} = [z^0, z^1_{l-1}, \dots, z^M_{l-1}]$$

  ② Interact with multi-head attention.
  $$[q, k, v] = \widehat{Z}_{l-1} W_{qkv}$$
  $$A = softmax(qk^T / C_h^{1/4})$$
  $$MSA(\widehat{Z}_{l-1}) = (Av)W_{out}$$

- Channel Interaction

  ① Average patch features: $z^c_{l-1} = \frac{1}{M} \sum_{i=1}^{M} z^i_{l-1}$

  ② Feed the prompt and visual context into MLP.
  $$\boldsymbol{\beta}_{l-1} = MLP([z^0; z^c_{l-1}])$$

  ③ Add the bias vector to all patch features.
  $$\widehat{Z}_{l-1} = [z^i_{l-1} + \boldsymbol{\beta}_{l-1},]  \quad i = 1, 2, \dots, M$$

# Experimental results

- miniImageNet & tieredImageNet

| Method | Backbone | Params/FLOPS | miniImageNet 5-way | | tieredImageNet 5-way | |
|---|---|---|---|---|---|---|
| | | | 1-shot | 5-shot | 1-shot | 5-shot |
| LEO [42] | WRN-28-10 | $36.5M/3.7 \times 10^{10}$ | $61.76\pm0.08$ | $77.59\pm0.12$ | $66.33\pm0.05$ | $81.44\pm0.09$ |
| CC+rot [14] | WRN-28-10 | $36.5M/3.7 \times 10^{10}$ | $62.93\pm0.45$ | $79.87\pm0.33$ | $70.53\pm0.51$ | $84.98\pm0.36$ |
| Align [1] | WRN-28-10 | $36.5M/3.7 \times 10^{10}$ | $65.92\pm0.60$ | $82.85\pm0.55$ | $\mathbf{74.40\pm0.68}$ | $86.61\pm0.59$ |
| MetaOptNet [22] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | $62.64\pm0.61$ | $78.63\pm0.46$ | $65.99\pm0.72$ | $81.56\pm0.53$ |
| Meta-Baseline [6] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | $63.17\pm0.23$ | $79.26\pm0.17$ | $68.62\pm0.27$ | $83.74\pm0.18$ |
| DeepEMD [56] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | $65.91\pm0.82$ | $82.41\pm0.56$ | $71.16\pm0.87$ | $86.03\pm0.58$ |
| RE-Net [17] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | $67.60\pm0.44$ | $82.58\pm0.30$ | $71.61\pm0.51$ | $85.28\pm0.35$ |
| TPMM [51] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | $67.64\pm0.63$ | $\mathbf{83.44\pm0.43}$ | $72.24\pm0.70$ | $86.55\pm0.63$ |
| SetFeat [2] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | $\mathbf{68.32\pm0.62}$ | $82.71\pm0.46$ | $73.63\pm0.88$ | $\mathbf{87.59\pm0.57}$ |
| SUN [10] | Visformer-S | $12.4M/1.7 \times 10^{8}$ | $67.80\pm0.45$ | $83.25\pm0.30$ | $72.99\pm0.50$ | $86.74\pm0.33$ |
| KTN [32] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | $61.42\pm0.72$ | $74.16\pm0.56$ | - | - |
| AM3 [52] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | $65.30\pm0.49$ | $78.10\pm0.36$ | $69.08\pm0.47$ | $82.58\pm0.31$ |
| TRAML [24] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | $\mathbf{67.10\pm0.52}$ | $79.54\pm0.60$ | - | - |
| DeepEMD-BERT [53] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | $67.03\pm0.79$ | $\mathbf{83.68\pm0.65}$ | $\mathbf{73.76\pm0.72}$ | $\mathbf{87.51\pm0.75}$ |
| Pre-train (Ours) | Visformer-T | $10.0M/1.3 \times 10^{9}$ | $65.16\pm0.44$ | $81.22\pm0.32$ | $72.38\pm0.50$ | $86.74\pm0.34$ |
| **SP-CLIP** (Ours) | Visformer-T | $10.0M/1.3 \times 10^{9}$ | $\mathbf{72.31\pm0.40}$ | $83.42\pm0.30$ | $\mathbf{78.03\pm0.46}$ | $88.55\pm0.32$ |
| **SP-SBERT** (Ours) | Visformer-T | $10.0M/1.3 \times 10^{9}$ | $70.70\pm0.42$ | $\mathbf{83.55\pm0.30}$ | $73.31\pm0.50$ | $88.56\pm0.32$ |
| **SP-GloVe** (Ours) | Visformer-T | $10.0M/1.3 \times 10^{9}$ | $70.81\pm0.42$ | $83.31\pm0.30$ | $74.68\pm0.50$ | $\mathbf{88.64\pm0.31}$ |

Table 1. Comparison with previous work on miniImageNet and tieredImageNet. Methods in the top rows do not use semantic information, and methods in the middle rows leverage semantic information from class names [24,32,52] or descriptions [53]. Accuracies are reported with 95% confidence intervals.

# Experimental results

- CIFAR-FS & FC100

| Method | Backbone | Params/FLOPs | CIFAR-FS 5-way | | FC100 5-way | |
|---|---|---|---|---|---|---|
| | | | 1-shot | 5-shot | 1-shot | 5-shot |
| PN+rot [14] | WRN-28-10 | $36.5M/3.7 \times 10^{10}$ | $69.55\pm0.34$ | $82.34\pm0.24$ | - | - |
| Align [1] | WRN-28-10 | $36.5M/3.7 \times 10^{10}$ | - | - | $\mathbf{45.83\pm0.48}$ | $59.74\pm0.56$ |
| ProtoNet [45] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | $72.2\pm0.7$ | $83.5\pm0.5$ | $37.5\pm0.6$ | $52.5\pm0.6$ |
| MetaOptNet [22] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | $72.6\pm0.7$ | $84.3\pm0.5$ | $41.1\pm0.6$ | $55.5\pm0.6$ |
| MABAS [18] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | $73.51\pm0.92$ | $85.49\pm0.68$ | $42.31\pm0.75$ | $57.56\pm0.78$ |
| Distill [47] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | $73.9\pm0.8$ | $86.9\pm0.5$ | $44.6\pm0.7$ | $\mathbf{60.9\pm0.6}$ |
| RE-Net [17] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | $74.51\pm0.46$ | $86.60\pm0.32$ | - | - |
| infoPatch [27] | ResNet-12 | $12.5M/3.5 \times 10^{9}$ | - | - | $43.8\pm0.4$ | $58.0\pm0.4$ |
| SUN [10] | Visformer-S | $12.4M/1.7 \times 10^{8}$ | $\mathbf{78.37\pm0.46}$ | $\mathbf{88.84\pm0.32}$ | - | - |
| Pre-train (Ours) | Visformer-T | $10.0M/1.3 \times 10^{9}$ | $71.99\pm0.47$ | $85.98\pm0.34$ | $43.77\pm0.39$ | $59.48\pm0.39$ |
| **SP-CLIP** (Ours) | Visformer-T | $10.0M/1.3 \times 10^{9}$ | $\mathbf{82.18\pm0.40}$ | $88.24\pm0.32$ | $\mathbf{48.53\pm0.38}$ | $\mathbf{61.55\pm0.41}$ |
| **SP-SBERT** (Ours) | Visformer-T | $10.0M/1.3 \times 10^{9}$ | $81.32\pm0.40$ | $88.31\pm0.32$ | $47.03\pm0.40$ | $61.03\pm0.40$ |
| **SP-GloVe** (Ours) | Visformer-T | $10.0M/1.3 \times 10^{9}$ | $81.62\pm0.41$ | $\mathbf{88.32\pm0.32}$ | $46.69\pm0.41$ | $61.18\pm0.41$ |

Table 2. Comparison with previous work on CIFAR-FS [22] and FC100 [31].

# Experimental results

| Aug | SI | CI | Mini | Tiered | CIFAR-FS | FC100 |
|-----|-----|-----|-------|--------|----------|-------|
| ✗ | ✗ | ✗ | 61.96 | 71.91 | 68.84 | 40.78 |
| ✓ | ✗ | ✗ | 65.15 | 72.38 | 71.99 | 43.77 |
| ✓ | ✓ | ✗ | 71.59 | 76.20 | 81.19 | 47.83 |
| ✓ | ✗ | ✓ | 70.48 | 77.62 | 79.80 | 47.10 |
| ✓ | ✓ | ✓ | **72.31** | **78.03** | **82.18** | **48.53** |

↑ 5.9%
↑ 5.4%
↑ 6.9%

Table 3. Ablation study on four datasets under the 1-shot setting. SI means spatial interaction, and CI means channel interaction.



Input image with harvestman and spider web | Pre-training baseline | Prompt with harvestman | Prompt with spider web
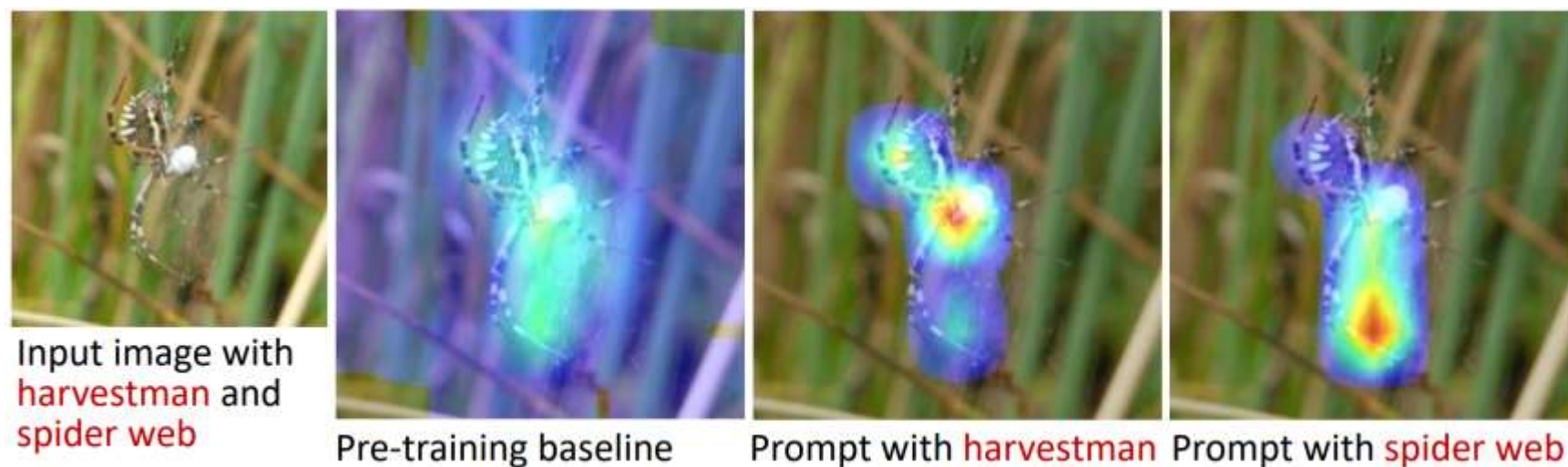
Figure 4. Visualization of attention maps when prompting with different class labels.
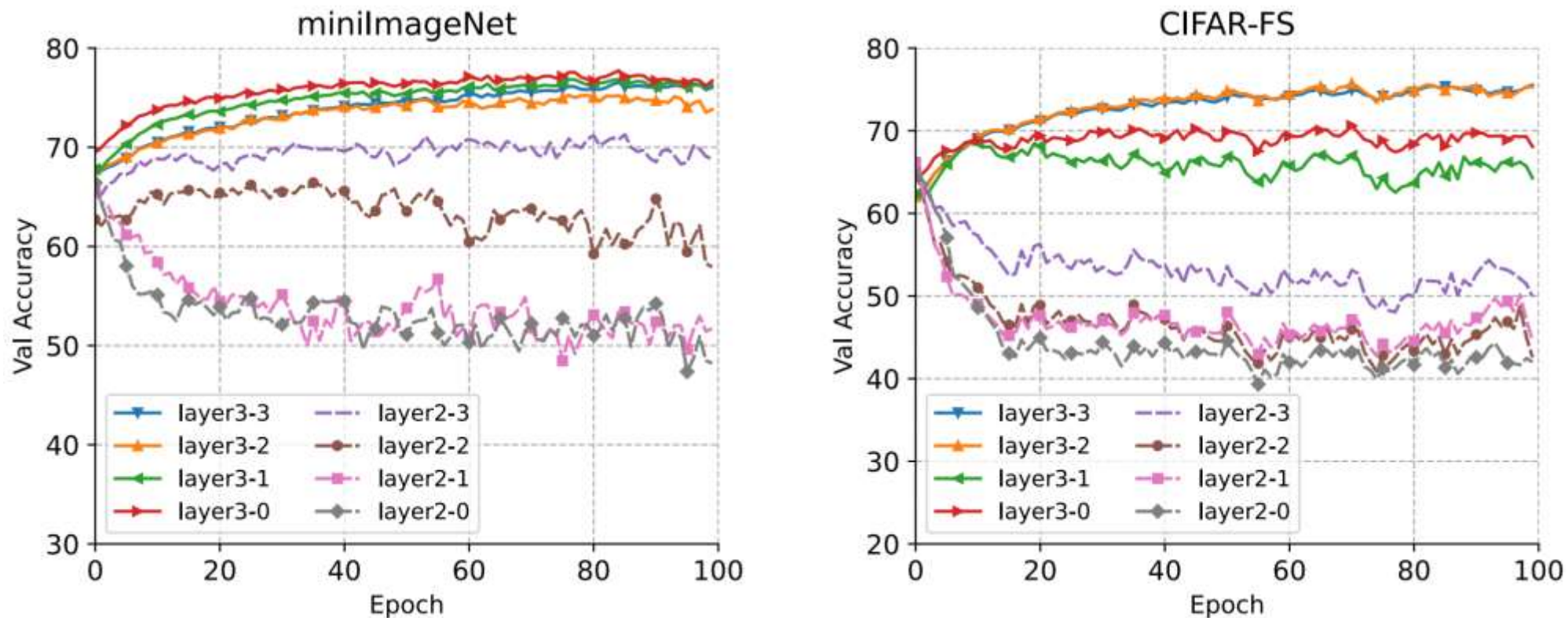
# Experimental results



Figure 3. Accuracy vs. different layers to inset prompts. We report 5-way 1-shot accuracy (%) on the validation set of miniImageNet and CIFAF-FS along the meta-training process. The feature extractor has three stages and multiple Transformer layers in each stage.
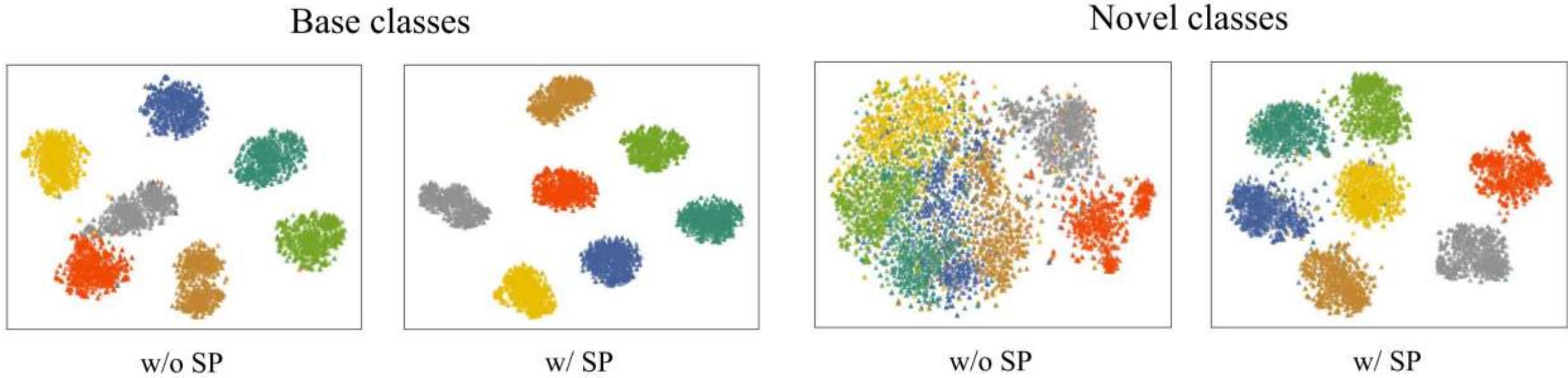
# Experimental results



Figure 5. t-SNE results of feature distributions.

# Summary

- We investigate how to use text data to improve the visual feature extraction for few-shot learning.

- We propose a new semantic prompt approach, where text features are used as prompts to adaptively tune the visual features.

- We propose two interaction mechanism, which allow the semantic prompt and visual features to interact along the spatial and the channel dimensions.

- Our approach is evaluated on four datasets with three different text encoders. Experimental results show that using semantic prompt can obtain much more performance gain than previous methods.