# CLIP-driven Coarse-to-fine Semantic Guidance for Fine-grained Open-set Semi-supervised Learning
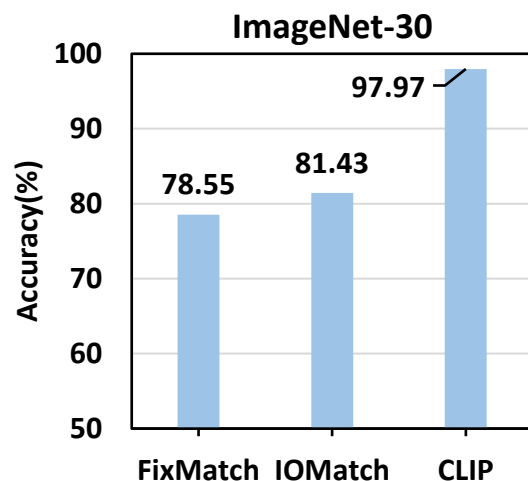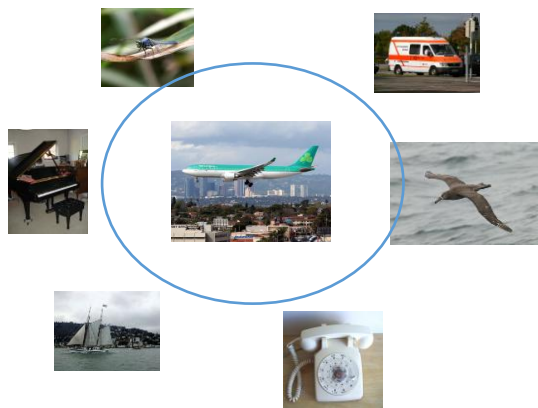
Xiaokun Li, Yaping Huang, Qingji Guan
Beijing Key Laboratory of Traffic Data Mining and Embodied Intelligence, Beijing Jiaotong University

Fine-grained open-set semi-supervised learning (OSSL) investigates a practical scenario where unlabeled data may contain fine-grained out-of-distribution (OOD) samples. Due to the subtle visual differences among in-distribution (ID) samples, as well as between ID and OOD samples, it is extremely challenging to separate the ID and OOD samples.
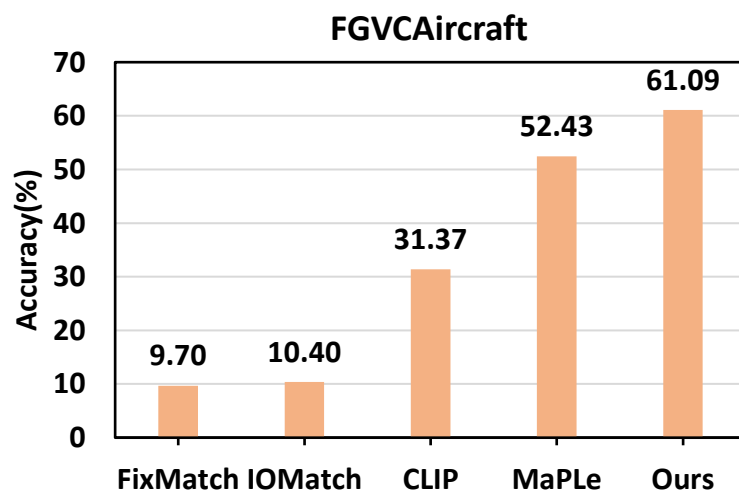
## Comparison of OSSL methods on coarse-grained and fine-grained classification tasks.

**(a) Coarse-grained OSSL**



ImageNet-30

- FixMatch: 78.55
- IOMatch: 81.43
- CLIP: 97.97

✓ On the coarse-grained ImageNet-30 dataset, some out-standing methods (e.g., FixMatch, IOMatch) achieve excellent performance. CLIP trained on the large-scale image-text pairs dataset achieves high performance as expected.

**(b) Fine-grained OSSL**



FGVCAircraft

- FixMatch: 9.70
- IOMatch: 10.40
- CLIP: 31.37
- MaPLe: 52.43
- Ours: 61.09

✓ On the fine-grained FGVCAircraft dataset, the generalization capabilities of CLIP remain effective but limited, as it tends to focus on general attributes while failing to capture the fine-grained details.

# Overview of the proposed CFSG-CLIP framework.



CFSG-CLIP is composed of a coarse-guidance branch and a fine-guidance branch based on the pre-trained CLIP model.

$t_c$ Coarse textual [CLS] tokens   $\tilde{z}_c$ Coarse visual [CLS] tokens   $z_c$ Coarse visual patch tokens   $\mathcal{S}$ Cosine Similarity

$t_f$ Fine textual [CLS] tokens   $\tilde{z}_f$ Fine visual [CLS] tokens   $z_f$ Fine visual patch tokens

## Semantic Filtering Module



**CLS** Image [CLS] token    $\mathcal{S}$ Similarity

$t_c$ Text [CLS] token    $w$ Weight & Sum

Top-k patch-level visual features:

$$s_{c_i} = \mathrm{sim}(z_{c_i}, t_c^m),$$

$$\mathcal{K} = \{i \in P : \mathrm{rank}(s_{c_i}) \leq k\},$$

Image-level global feature weighting:
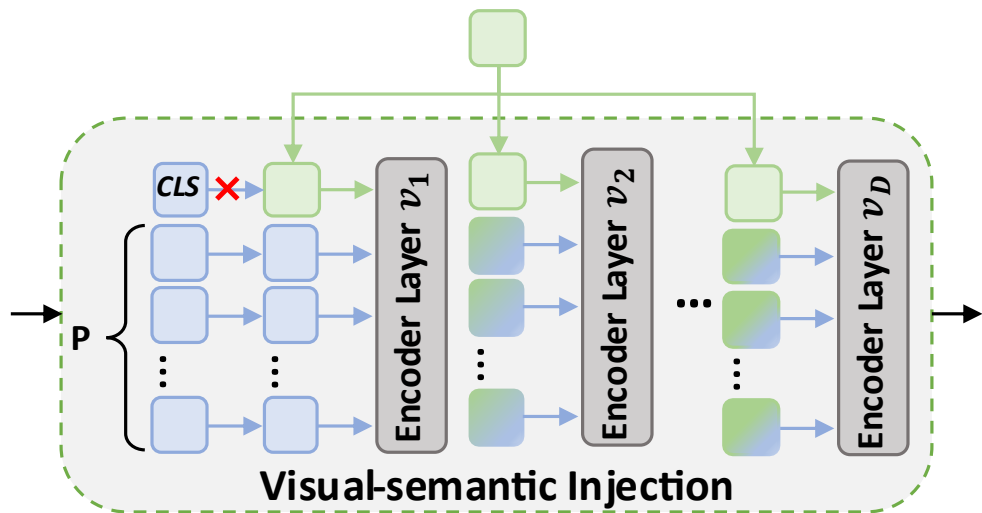
$$w_{c_i} = \frac{\exp(\mathrm{sim}(z_{c_i}, \tilde{z}_c))}{\sum \exp(\mathrm{sim}(z_{c_i}, \tilde{z}_c))}, i \in \mathcal{K}$$

$$z_c = \sum_{i \in \mathcal{K}} w_{c_i} z_{c_i}$$

To capture fine-grained details and maintain semantic consistency, we design a semantic filtering module that leverages both textual and visual features to obtain global and local visual features.

5

## Semantic Filtering Module



**Visual-semantic Injection**

Injecting the local visual features into transformer blocks:

$$[\_, E_j] = \mathcal{V}_j([\text{proj}(z_c), E_{j-1}]) \quad j = 1, 2, \ldots, D,$$

$$[\tilde{z}_{f_j}, E_j] = \mathcal{V}_j([\text{proj}(z_c), E_{j-1}]) \quad j = D+1, \ldots, T,$$

To guide the visual encoder to focus more on the fine-grained details of the image, we further design a visual semantic injection strategy in the fine-guidance branch.

## Dual-branch Training

Optimization of Coarse-guidance Branch:

$$\tilde{p}_c^l = \frac{\exp(\mathrm{sim}(\tilde{z}_c^l, t_c^m)/\tau)}{\sum_{m'} \exp(\mathrm{sim}(\tilde{z}_c^l, t_c^{m'})/\tau)},$$

$$p_c^l = \frac{\exp(\mathrm{sim}(z_c^l, t_c^m)/\tau)}{\sum_{m'} \exp(\mathrm{sim}(z_c^l, t_c^{m'})/\tau)},$$

$$L_c = H(y^l, \tilde{p}_c^l) + H(y^l, p_c^l) + \lambda_c(\mathcal{F}(x^u)H(\tilde{p}_c^{u_w}, \tilde{p}_c^{u_s})$$
$$+ \mathcal{F}(x^u)H(p_c^{u_w}, p_c^{u_s})),$$

## Dual-branch Training

Optimization of Fine-guidance Branch:

$$\tilde{p}_f^l = \frac{\exp(\mathrm{sim}(\tilde{z}_f^l, t_f^m)/\tau)}{\sum_{m'} \exp(\mathrm{sim}(\tilde{z}_f^l, t_f^{m'})/\tau)}.$$

$$p_f^l = \frac{\exp(\mathrm{sim}(z_f^l, t_f^m)/\tau)}{\sum_{m'} \exp(\mathrm{sim}(z_f^l, t_f^{m'})/\tau)}.$$

$$L_f = H(y^l, \tilde{p}_f^l) + H(y^l, p_f^l) + \lambda_f(\mathcal{F}(\tilde{u}_f)H(\tilde{p}_f^{u_w}, \tilde{p}_f^{u_s})$$
$$+ \mathcal{F}(u_f)H(p_f^{u_w}, p_f^{u_s})),$$

## Dual-branch Training

Optimization of Fine-guidance Branch:

$$\tilde{p}_f^l = \frac{\exp(\mathrm{sim}(\tilde{z}_f^l, t_f^m)/\tau)}{\sum_{m'} \exp(\mathrm{sim}(\tilde{z}_f^l, t_f^{m'})/\tau)}.$$

$$p_f^l = \frac{\exp(\mathrm{sim}(z_f^l, t_f^m)/\tau)}{\sum_{m'} \exp(\mathrm{sim}(z_f^l, t_f^{m'})/\tau)}.$$

$$L_f = H(y^l, \tilde{p}_f^l) + H(y^l, p_f^l) + \lambda_f(\mathcal{F}(\tilde{u}_f)H(\tilde{p}_f^{u_w}, \tilde{p}_f^{u_s})$$
$$+ \mathcal{F}(u_f)H(p_f^{u_w}, p_f^{u_s})),$$

Table 1. Classification accuracy (%) for CLIP-based methods on four fine-grained benchmark datasets with varying labeled set sizes under the fine-grained OSSL setting. The results are presented as the mean with standard deviation over three runs using different random seeds.

| Method | Stanford Dogs | | Stanford Cars | | CUB-200-2011 | | FGVCAircraft | |
|---|---|---|---|---|---|---|---|---|
| | 5 | 20 | 5 | 20 | 5 | 20 | 5 | 20 |
| CLIP [27] | $79.25_{\pm 0.00}$ | $79.25_{\pm 0.00}$ | $75.97_{\pm 0.00}$ | $75.97_{\pm 0.00}$ | $66.00_{\pm 0.00}$ | $66.00_{\pm 0.00}$ | $31.37_{\pm 0.00}$ | $31.37_{\pm 0.00}$ |
| CLIP-LORA [41] | $83.81_{\pm 0.37}$ | $84.31_{\pm 0.27}$ | $82.71_{\pm 0.56}$ | $82.45_{\pm 1.30}$ | $70.95_{\pm 0.85}$ | $73.40_{\pm 0.75}$ | $40.89_{\pm 1.73}$ | $42.37_{\pm 0.71}$ |
| CLIP-Adapter [7] | $82.91_{\pm 0.25}$ | $86.02_{\pm 0.27}$ | $84.31_{\pm 0.02}$ | $87.13_{\pm 0.28}$ | $80.03_{\pm 0.29}$ | $84.77_{\pm 1.02}$ | $47.77_{\pm 0.90}$ | $55.79_{\pm 0.73}$ |
| CoOp [45] | $83.01_{\pm 0.26}$ | $85.68_{\pm 0.37}$ | $85.45_{\pm 0.31}$ | $87.64_{\pm 0.46}$ | $80.10_{\pm 0.29}$ | $85.40_{\pm 0.37}$ | $45.39_{\pm 0.96}$ | $55.43_{\pm 0.30}$ |
| LoCoOp [24] | $83.08_{\pm 0.25}$ | $86.26_{\pm 0.11}$ | $84.10_{\pm 0.72}$ | $87.83_{\pm 0.66}$ | $79.27_{\pm 0.45}$ | $85.63_{\pm 0.54}$ | $45.53_{\pm 1.36}$ | $54.67_{\pm 1.59}$ |
| PLOT [3] | $84.46_{\pm 0.07}$ | $87.11_{\pm 0.09}$ | $86.28_{\pm 0.30}$ | $88.59_{\pm 0.45}$ | $81.43_{\pm 0.66}$ | $87.20_{\pm 0.14}$ | $49.59_{\pm 0.37}$ | $58.25_{\pm 0.93}$ |
| MaPLe [14] | $\mathbf{85.64}_{\pm 0.15}$ | $87.64_{\pm 0.20}$ | $88.16_{\pm 0.25}$ | $90.34_{\pm 0.25}$ | $83.30_{\pm 0.33}$ | $88.77_{\pm 0.21}$ | $52.43_{\pm 0.47}$ | $64.33_{\pm 1.21}$ |
| Ours | $85.48_{\pm 0.21}$ | $\mathbf{89.42}_{\pm 0.16}$ | $\mathbf{90.38}_{\pm 0.09}$ | $\mathbf{93.08}_{\pm 0.08}$ | $\mathbf{84.73}_{\pm 0.17}$ | $\mathbf{91.75}_{\pm 0.24}$ | $\mathbf{61.09}_{\pm 0.27}$ | $\mathbf{73.56}_{\pm 0.58}$ |

Table 2. Open-set classification balanced accuracy (%) for CLIP-based methods on four fine-grained benchmark datasets with varying labeled set sizes under the fine-grained OSSL setting.

| Method | Stanford Dogs | | Stanford Cars | | CUB-200-2011 | | FGVCAircraft | |
|---|---|---|---|---|---|---|---|---|
| | 5 | 20 | 5 | 20 | 5 | 20 | 5 | 20 |
| CLIP [27] | $77.17_{\pm 0.00}$ | $77.17_{\pm 0.00}$ | $75.70_{\pm 0.00}$ | $75.70_{\pm 0.00}$ | $64.10_{\pm 0.00}$ | $64.10_{\pm 0.00}$ | $31.08_{\pm 0.00}$ | $31.08_{\pm 0.00}$ |
| CLIP-LORA [41] | $82.34_{\pm 0.57}$ | $82.67_{\pm 0.36}$ | $82.10_{\pm 0.56}$ | $81.03_{\pm 1.24}$ | $70.52_{\pm 1.34}$ | $72.20_{\pm 0.45}$ | $40.10_{\pm 1.69}$ | $41.57_{\pm 0.69}$ |
| CLIP-Adapter [7] | $81.63_{\pm 0.14}$ | $84.36_{\pm 0.25}$ | $83.65_{\pm 0.04}$ | $86.50_{\pm 0.25}$ | $81.36_{\pm 0.60}$ | $85.20_{\pm 0.95}$ | $46.87_{\pm 0.88}$ | $53.52_{\pm 1.35}$ |
| CoOp [45] | $81.65_{\pm 0.20}$ | $84.25_{\pm 0.41}$ | $84.63_{\pm 0.36}$ | $86.84_{\pm 0.42}$ | $81.04_{\pm 0.06}$ | $84.92_{\pm 0.51}$ | $44.55_{\pm 0.94}$ | $54.35_{\pm 0.30}$ |
| LoCoOp [24] | $81.78_{\pm 0.23}$ | $84.68_{\pm 0.20}$ | $83.25_{\pm 0.76}$ | $87.17_{\pm 0.64}$ | $80.45_{\pm 0.86}$ | $85.46_{\pm 0.32}$ | $44.67_{\pm 1.35}$ | $53.60_{\pm 1.57}$ |
| PLOT [3] | $82.95_{\pm 0.07}$ | $85.54_{\pm 0.11}$ | $85.58_{\pm 0.32}$ | $87.83_{\pm 0.34}$ | $83.32_{\pm 0.31}$ | $87.16_{\pm 0.17}$ | $48.68_{\pm 0.37}$ | $57.13_{\pm 0.92}$ |
| MaPLe [14] | $\mathbf{84.09}_{\pm 0.19}$ | $86.02_{\pm 0.25}$ | $87.43_{\pm 0.27}$ | $89.48_{\pm 0.25}$ | $84.72_{\pm 0.53}$ | $88.66_{\pm 0.60}$ | $51.79_{\pm 0.43}$ | $63.12_{\pm 1.19}$ |
| Ours | $84.02_{\pm 0.15}$ | $\mathbf{87.77}_{\pm 0.19}$ | $\mathbf{89.65}_{\pm 0.05}$ | $\mathbf{92.34}_{\pm 0.10}$ | $\mathbf{86.46}_{\pm 0.25}$ | $\mathbf{90.92}_{\pm 0.24}$ | $\mathbf{59.92}_{\pm 0.26}$ | $\mathbf{72.13}_{\pm 0.57}$ |

Table 3. Classification accuracy (%) on the Semi-Aves dataset is reported for two settings: unlabeled data with ID samples $U_{in}$, and unlabeled data with a mix of ID and OOD samples $U_{in} + U_{out}$.

| Method | $U_{in}$ | $U_{in} + U_{out}$ |
|---|---|---|
| CLIP [27] | $10.05 \pm 0.00$ | $10.05 \pm 0.00$ |
| CLIP-Adapter [7] | $50.16 \pm 1.55$ | $50.10 \pm 0.44$ |
| CoOp [45] | $47.67 \pm 0.97$ | $47.93 \pm 0.76$ |
| LoCoOp [24] | $48.04 \pm 0.96$ | $47.30 \pm 1.30$ |
| PLOT [3] | $53.68 \pm 0.35$ | $54.72 \pm 0.96$ |
| MaPLe [14] | $60.39 \pm 0.10$ | $59.68 \pm 0.64$ |
| Ours | $\mathbf{65.63} \pm \mathbf{0.27}$ | $\mathbf{63.69} \pm \mathbf{0.74}$ |

Table 4. Ablation studies on the Stanford cars and FGVCAircraft datasets. The 'A' stands for adapter, 'SFM' denotes semantic filtering module, 'C' means coarse-guidance branch, 'F' means fine-guidance branch, and 'VSI' is visual-semantic injection strategy. The results are reported based on a single run with seed 1.

| Method | Stanford Cars | FGVCAircraft |
|---|---|---|
| CoOp+$\mathcal{A}$ | 86.31 | 50.42 |
| +SFM (C) | 89.15 | 54.65 |
| +SFM (C)+SFM (F) | 90.03 | 58.01 |
| +SFM (C)+SFM (F)+VSI | 90.28 | 61.07 |

Table 5. Ablation studies for semantic filtering module. The results are reported based on a single run with seed 1.

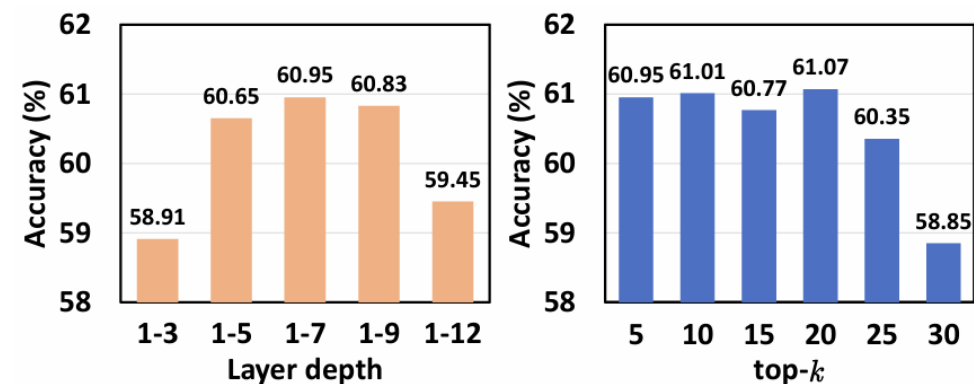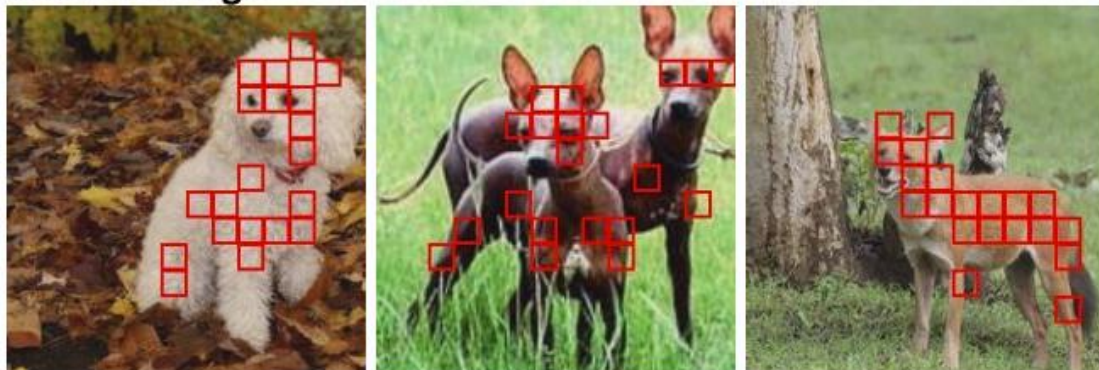| Method | Stanford Cars | FGVCAircraft |
|---|---|---|
| CoOp+$\mathcal{A}$ | 86.31 | 50.42 |
| Textual Filtering | 88.65 | 53.57 |
| Visual Weighting | 89.15 | 54.65 |



Figure 6. Ablation studies on injection depth (*left*) and top-$k$ (*right*) on the FGVCAircraf datasets. The results are reported based on a single run with seed 1.

Table 6. Evaluation results for different operations of employing local visual features. The results are reported based on a single run with seed 1.

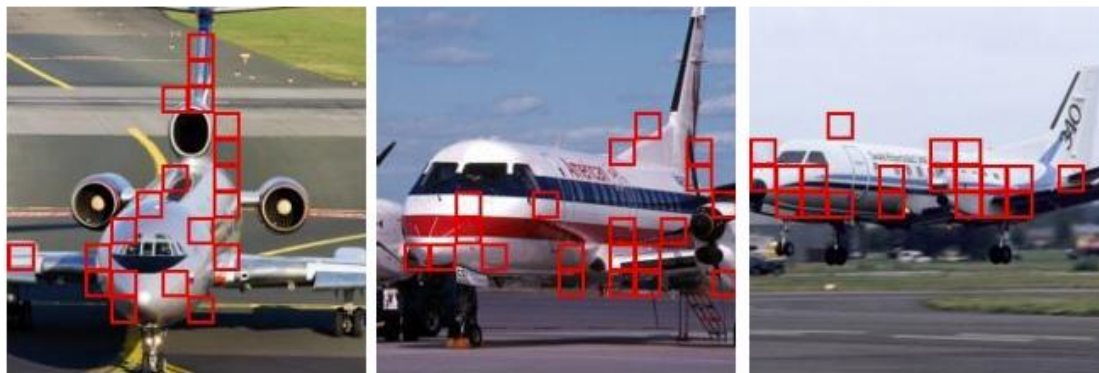| Method | Stanford Cars | FGVCAircraft |
|---|---|---|
| Concatenate | 89.99 | 60.47 |
| Replace | 90.28 | 61.07 |

Figure 5. Visualization of patch-tokens extracted by semantic filtering module. We find that the semantic filtering module can correctly extract local visual regions on different fine-grained datasets.

# Thank you!