

# Exploring Structured Semantic Prior for Multi Label Recognition with Incomplete Labels

Zixuan Ding<sup>1,4\*</sup> Ao Wang<sup>2,3,4\*</sup> Hui Chen<sup>2,3,†</sup> Qiang Zhang<sup>1</sup>  
Pengzhang Liu<sup>5</sup> Yongjun Bao<sup>5</sup> Weipeng Yan<sup>5</sup> Jungong Han<sup>6,7</sup>

<sup>1</sup>Xidian University <sup>2</sup>Tsinghua University <sup>3</sup>BNRist

<sup>4</sup>Hangzhou Zhuoxi Institute of Brain and Intelligence <sup>5</sup>JD.com

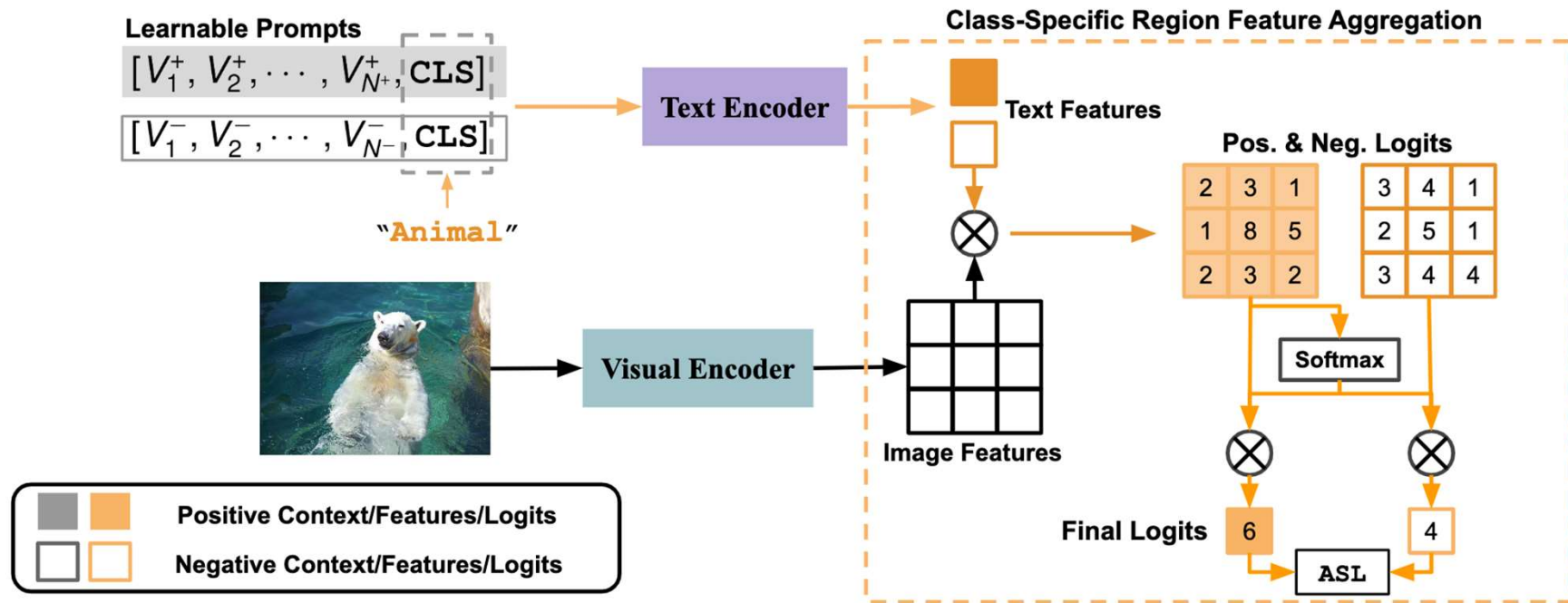
<sup>6</sup>Department of Computer Science, the University of Sheffield, UK

<sup>7</sup>Centre for Machine Intelligence, the University of Sheffield, UK

dingzixuan@stu.xidian.edu.cn wa22@mails.tsinghua.edu.cn qzhang@xidian.edu.cn  
{jichenhui2012, jungonghan77}@gmail.com {Paul.yan, baoyongjun, liupengzhang}@jd.com

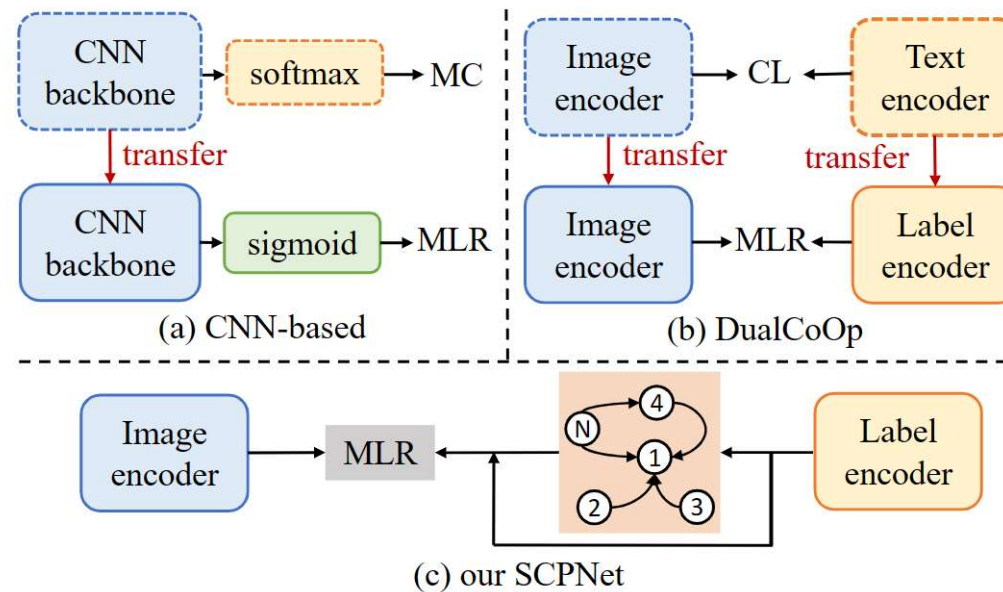
# Background

DualCoOp learns a pair of positive and negative prompts to quickly adapt powerful pretrained vision-text encoders to the MLR task. For each class, two prompts generate two contrastive (positive and negative) textual embeddings as the input to the text encoder. During training, we apply the ASL loss to optimize learnable prompts while keeping other network components frozen. During inference, we compare the positive and negative logits to make a prediction for each class

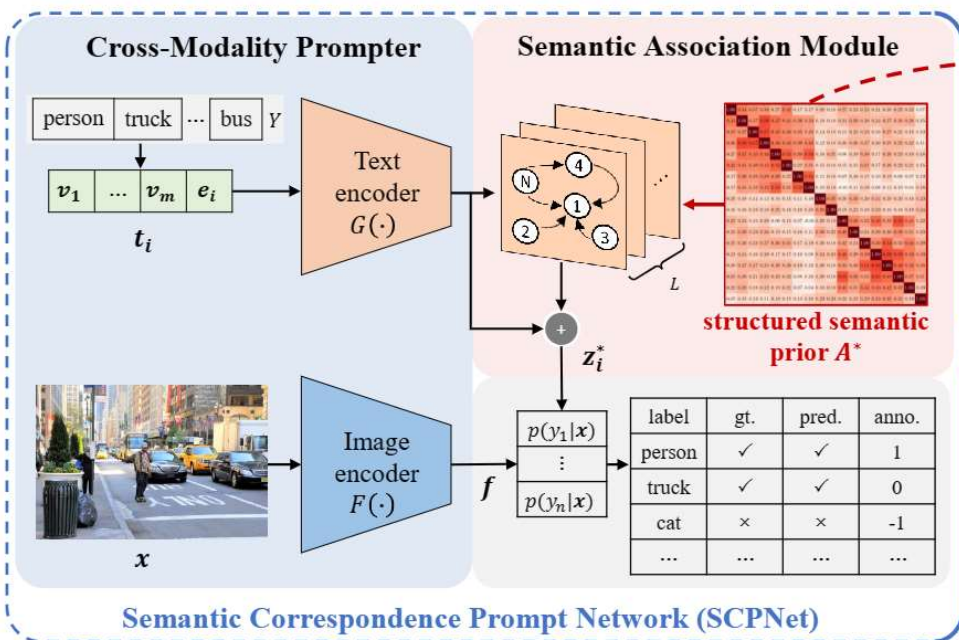


# Background

Overview of CNN-based, DualCoOp and our SCPNet. Like DualCoOp, our SCPNet adopts CLIP as the base model. Differently, our SCPNet aims to enhance the MLR with the prior about the label-to-label correspondence. *MC* means multi-class. *CL* denotes contrastive learning.



# Methods



## Cross-modality prompter (CMP):

Prompt:  $t_i = \{v_1, v_2, \dots, v_m, e_i\}$

soft prompt token

embedding of  $y_i$

$$f = F(x), z_i = G(t_i)$$

## Structured Prior Prompter

label feature

$$a_{ij} = \text{sim}(\bar{z}_i, \bar{z}_j)$$

mitigate the over-smoothness of graph representation by adjusting the sparse graph:

$$a'_{ij} = \begin{cases} a_{ij}, & \text{if } j \in \text{topK}(a_i) \\ 0, & \text{if } j \notin \text{topK}(a_i) \end{cases}$$

hyper-parameter

$$\bar{a}_{ij} = \begin{cases} (s / \sum_{i \neq j'} a'_{ij'}) \times a'_{ij}, & \text{if } i \neq j \\ 1 - s, & \text{if } i = j \end{cases}$$

$$a^*_{ij} = \frac{\mathbb{I}[\bar{a}_{ij} \neq 0] \exp(\bar{a}_{ij} / \tau')}{\sum_j \mathbb{I}[\bar{a}_{ij} \neq 0] \exp(\bar{a}_{ij} / \tau')}$$

## Semantic association module (SAM)

$$\mathbf{Z} = \{z_0, z_1, \dots, z_n\}$$

$$\text{input features } H^0 = \mathbf{Z}$$

a learnable parameter matrix

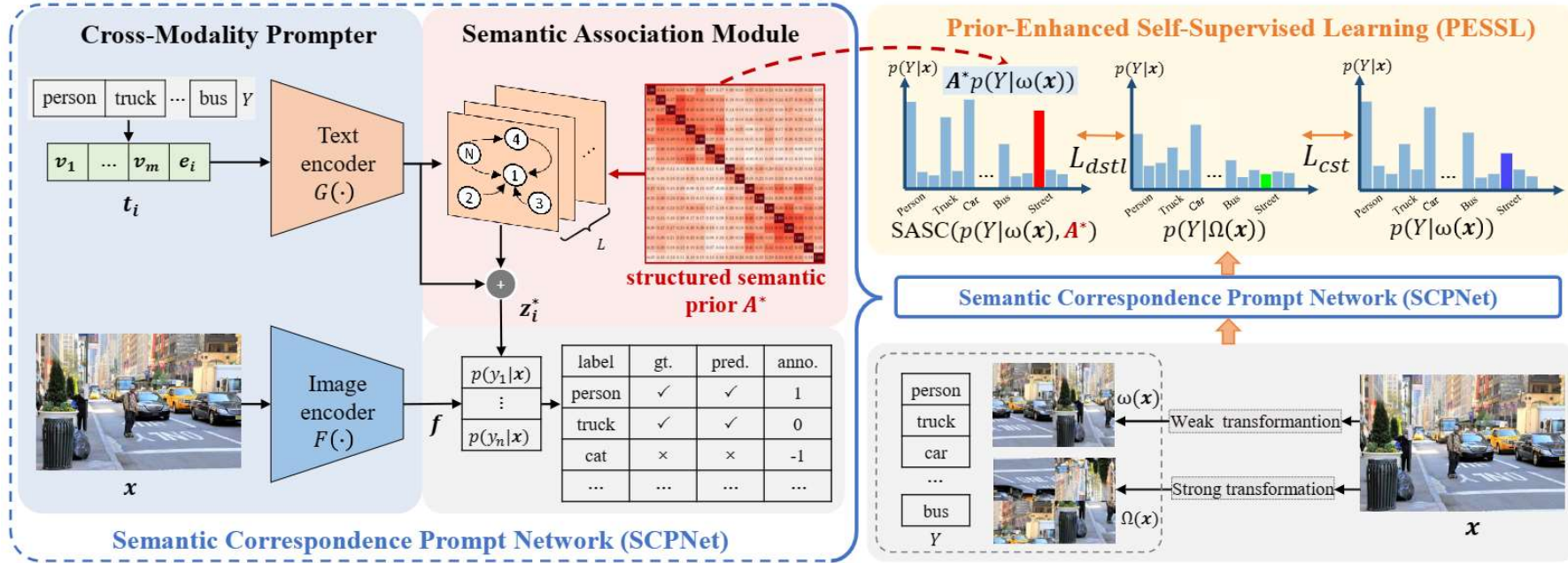
$$\text{GCN: } H^{l+1} = \rho(A^* H^l W^l)$$

non-linear function

$$\mathbf{Z}^* = H^0 + H^L$$

$$p(y_i | x) = \sigma(\text{sim}(f, z_i^*) / \tau)$$

# Methods



## Prior-Enhanced Self-Supervised Learning

$$p^*(y_i|x) = \sum_{y_j \in \mathcal{N}(y_i)} w(i, j) \times p(y_j|x)$$

a correlated neighboring set

$$\mathcal{L}_{cst} = - \sum_{c \in \mathcal{O}(x)} \log p(c|\Omega(x)) - \sum_{c \notin \mathcal{O}(x)} \log(1 - p(c|\Omega(x)))$$

$$\mathcal{L}_{pessl} = \lambda_{cst} \mathcal{L}_{cst} + \lambda_{dstl} \mathcal{L}_{dstl}$$

$$p^*(y|\omega(x)) = \text{SASC}(p(y|\omega(x)), A^*)$$

$$\mathcal{L} = \mathcal{L}_{cls} + \mathcal{L}_{pessl}$$

We then customize the whole process as a function parameterized by W:

$$\text{SASC}(p(Y|x), \mathbb{W}) = \mathbb{W}p(Y|x)$$

$$\mathcal{L}_{dstl} = - \sum_c \left( q_c^w \log \frac{q_c^s}{q_c^w} + (1 - q_c^w) \log \frac{1 - q_c^s}{1 - q_c^w} \right)$$

where  $q_c^w = p^*(c|\omega(x))$  and  $q_c^s = p(c|\Omega(x))$

# Experiments



Method	LargeLoss setup [16]					SPLC setup [32]				
	COCO	VOC	NUS	CUB	Avg.	COCO	VOC	NUS	CUB	Avg.
LSAN [8]	69.2	86.7	50.5	17.9	56.1	70.5	87.2	52.5	18.9	57.3
ROLE [8]	69.0	88.2	51.0	16.8	56.3	70.9	89.0	50.6	20.4	57.7
LargeLoss [16]	71.6	89.3	49.6	21.8	58.1	-	-	-	-	-
Hill [32]	-	-	-	-	-	73.2	87.8	55.0	18.8	58.7
SPLC [32]	72.0	87.7	49.8	18.0	56.9	73.2	88.1	55.2	20.0	59.1
<b>SCPNet (ours)</b>	<b>75.4</b>	<b>90.1</b>	<b>55.7</b>	<b>25.4</b>	<b>61.7</b>	<b>76.4</b>	<b>91.2</b>	<b>62.0</b>	<b>25.7</b>	<b>63.8</b>

single positive label

Datasets	Method	10%	20%	30%	40%	50%	60%	70%	80%	90%	Avg.
COCO	SSGRL [5]	62.5	70.5	73.2	74.5	76.3	76.5	77.1	77.9	78.4	74.1
	GCN-ML [6]	63.8	70.9	72.8	74.0	76.7	77.1	77.3	78.3	78.6	74.4
	SST [4]	68.1	73.5	75.9	77.3	78.1	78.9	79.2	79.6	79.9	76.7
	SARB [21]	71.2	75.0	77.1	78.3	78.9	79.6	79.8	80.5	80.5	77.9
	DualCoOp [26]	78.7	80.9	81.7	82.0	82.5	82.7	82.8	83.0	83.1	81.9
	<b>SCPNet (ours)*</b>	<b>80.3</b>	<b>82.2</b>	<b>82.8</b>	83.4	83.8	83.9	84.0	84.1	84.2	83.2
	<b>SCPNet (ours)</b>	79.1	82.1	<b>82.8</b>	<b>83.9</b>	<b>84.5</b>	<b>84.9</b>	<b>85.4</b>	<b>85.7</b>	<b>85.9</b>	<b>83.8</b>
VOC2007	SSGRL [5]	77.7	87.6	89.9	90.7	91.4	91.8	91.9	92.2	92.2	89.5
	GCN-ML [6]	74.5	87.4	89.7	90.7	91.0	91.3	91.5	91.8	92.0	88.9
	SST [4]	81.5	89.0	90.3	91.0	91.6	92.0	92.5	92.6	92.7	90.4
	SARB [21]	83.5	88.6	90.7	91.4	91.9	92.2	92.6	92.8	92.9	90.7
	DualCoOp [26]	90.3	92.2	92.8	93.3	93.6	93.9	94.0	94.1	94.2	93.2
	<b>SCPNet (ours)</b>	<b>91.1</b>	<b>92.8</b>	<b>93.5</b>	<b>93.6</b>	<b>93.8</b>	<b>94.0</b>	<b>94.1</b>	<b>94.2</b>	<b>94.3</b>	<b>93.5</b>
	<b>SCPNet (ours)</b>	<b>91.1</b>	<b>92.8</b>	<b>93.5</b>	<b>93.6</b>	<b>93.8</b>	<b>94.0</b>	<b>94.1</b>	<b>94.2</b>	<b>94.3</b>	<b>93.5</b>
VG-200	SSGRL [5]	34.6	37.3	39.2	40.1	40.4	41.0	41.3	41.6	42.1	39.7
	GCN-ML [6]	32.0	37.8	38.8	39.1	39.6	40.0	41.9	42.3	42.5	39.3
	SST [4]	38.8	39.4	41.1	41.8	42.7	42.9	43.0	43.2	43.5	41.8
	SARB [21]	41.4	44.0	44.8	45.5	46.6	47.5	47.8	48.0	48.2	46.0
	<b>SCPNet (ours)</b>	<b>43.8</b>	<b>46.4</b>	<b>48.2</b>	<b>49.6</b>	<b>50.4</b>	<b>50.9</b>	<b>51.3</b>	<b>51.6</b>	<b>52.0</b>	<b>49.4</b>
	<b>SCPNet (ours)</b>	<b>43.8</b>	<b>46.4</b>	<b>48.2</b>	<b>49.6</b>	<b>50.4</b>	<b>50.9</b>	<b>51.3</b>	<b>51.6</b>	<b>52.0</b>	<b>49.4</b>

MLR with partial labels

# Experiments



different modules in the proposed SCPNet method for both the single positive label setting and the partial label setting.

Model	CMP	SAM	PESSL		Single Positive Label				Partial Label			Avg.
			$\mathcal{L}_{cst}$	$\mathcal{L}_{dstl}$	COCO	VOC	NUS	CUB	COCO	VOC2007	VG-200	
Baseline					73.18	88.07	55.18	19.99	77.41	88.32	46.39	64.08
SCPNet	✓				74.36	88.46	60.66	21.42	80.90	89.16	47.55	66.07
	✓	✓			75.12	89.09	61.08	21.66	82.12	90.16	48.11	66.76
	✓	✓	✓		75.70	90.92	61.75	23.67	82.85	92.50	48.70	68.01
	✓	✓		✓	75.84	90.92	61.56	24.51	83.35	93.21	48.83	68.32
	✓	✓	✓	✓	<b>76.42</b>	<b>91.16</b>	<b>62.04</b>	<b>25.71</b>	<b>83.76</b>	<b>93.49</b>	<b>49.36</b>	<b>68.85</b>

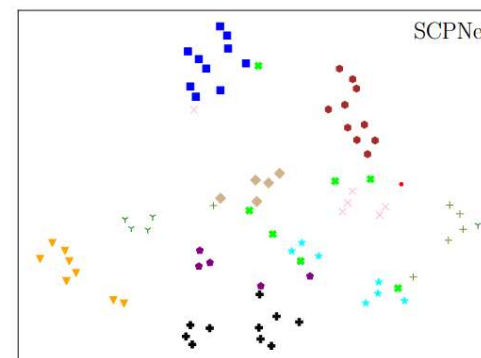
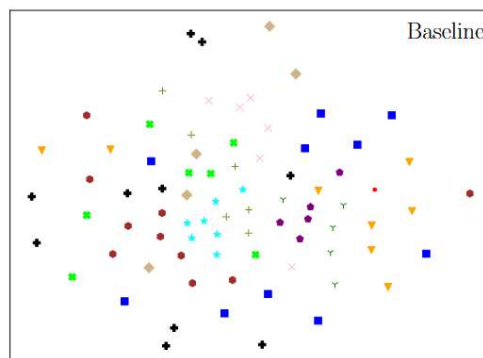
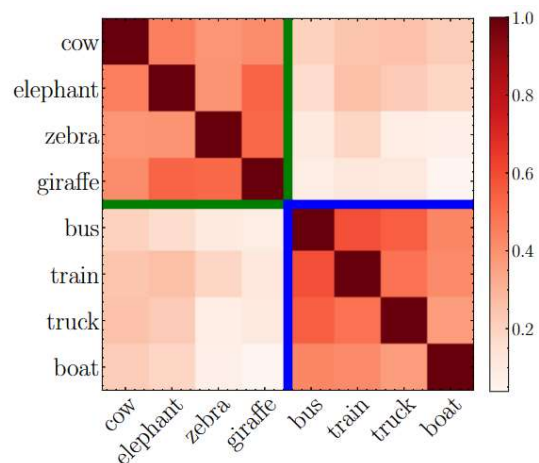
$$\mathcal{L}_{pessl} = \lambda_{cst}\mathcal{L}_{cst} + \lambda_{dstl}\mathcal{L}_{dstl}$$

$$\mathcal{L} = \mathcal{L}_{cls} + \mathcal{L}_{pessl}$$

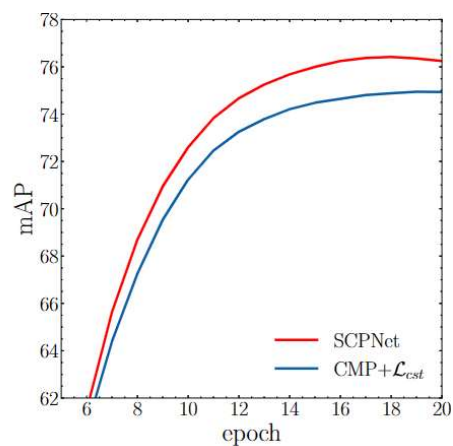
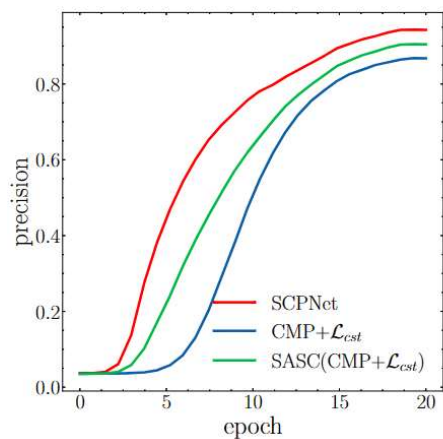
$L$	2	3	4
mAP	75.88	<b>76.42</b>	76.22

$\lambda_{cst}$	0	1/16	1/8	1/4	1/8		
$\lambda_{dstl}$	0				1/8	2/8	3/8
mAP	75.12	75.56	<b>75.70</b>	75.13	<b>76.42</b>	76.40	76.17

# Experiments



- person
- ▽ vehicle
- animal
- outdoor
- accessory
- kitchen
- furniture
- appliance
- sports
- food
- × electronic
- indoor



$CMP+L_{cst}$ : wipes out components involving the prior

It indicates that the quality of label supervision can be promoted under the guidance of the proposed prior, thus benefiting the performance on the test set.

Thanks