



# Multi-Label Image Recognition with Graph Convolutional Networks

---

Zhao-Min Chen<sup>1,2</sup>      Xiu-Shen Wei<sup>2</sup>      Peng Wang<sup>3</sup>      Yanwen Guo<sup>1</sup>

<sup>1</sup>National Key Laboratory for Novel Software Technology, Nanjing University, China

<sup>2</sup>Megvii Research Nanjing, Megvii Technology, China

<sup>3</sup>School of Computer Science, The University of Adelaide, Australia

---

# Introduction

## Multi-Label Learning



Person, Sports Ball,  
Tennis Racket



Person, Tie



Person, Ski

- Label Dependency

Label co-occurrence is often used to model the dependency.

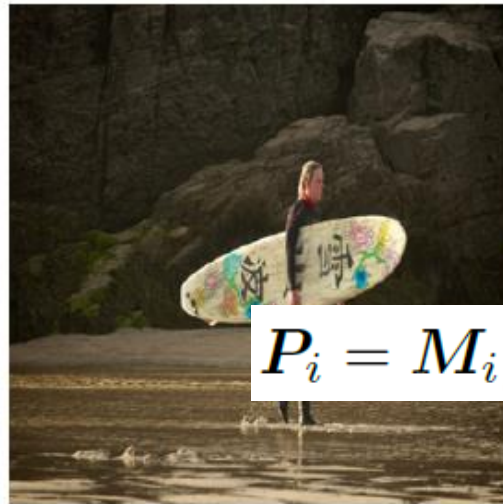
# Method

**Q1:** How to effectively capture the correlation between object labels?

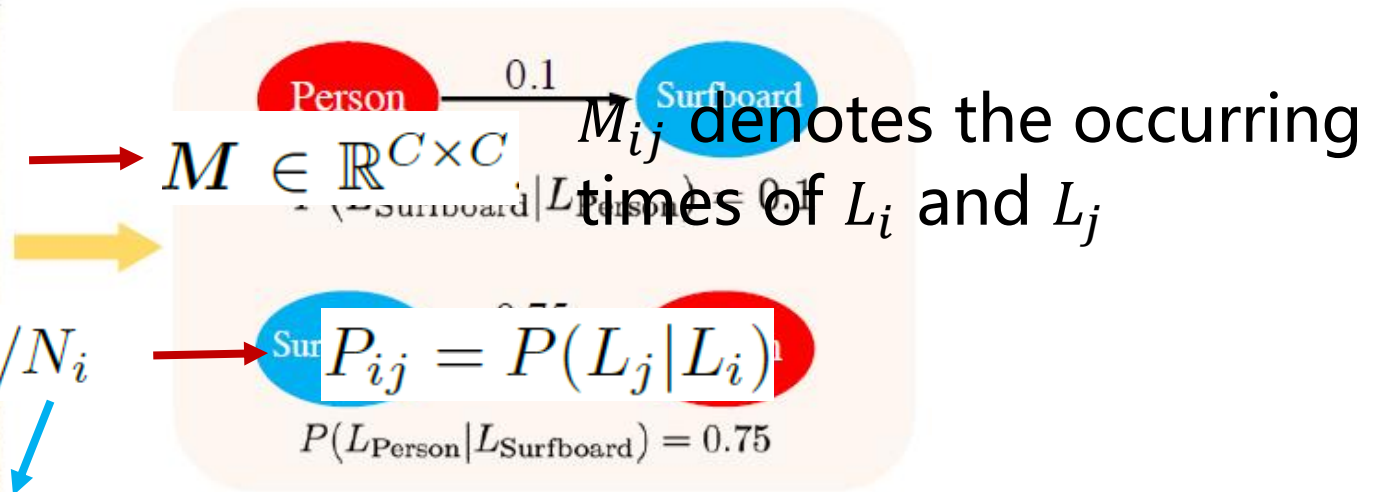
$P(L_j|L_i)$  → The probability of occurrence of label  $L_j$  when label  $L_i$  appears

## Correlation Matrix

count the occurrence of label pairs in the training set



The occurrence times of  $L_i$  in the training set



# Method

**Q1:** How to effectively capture the correlation between object labels?

Binarize the correlation

$$A_{ij} = \begin{cases} 0, & \text{if } P_{ij} < \tau \\ 1, & \text{if } P_{ij} \geq \tau \end{cases}$$

Re-weighted Scheme

$$A'_{ij} = \begin{cases} p / \sum_{\substack{j=1 \\ i \neq j}}^C A_{ij}, & \text{if } i \neq j \\ 1 - p, & \text{if } i = j \end{cases}$$



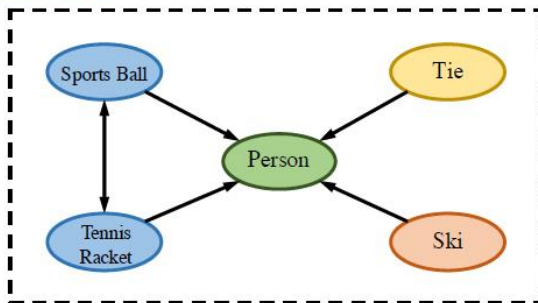
Person, Sports Ball,  
Tennis Racket

Person, Tie

Person, Ski

Label<sub>A</sub> → Label<sub>B</sub>

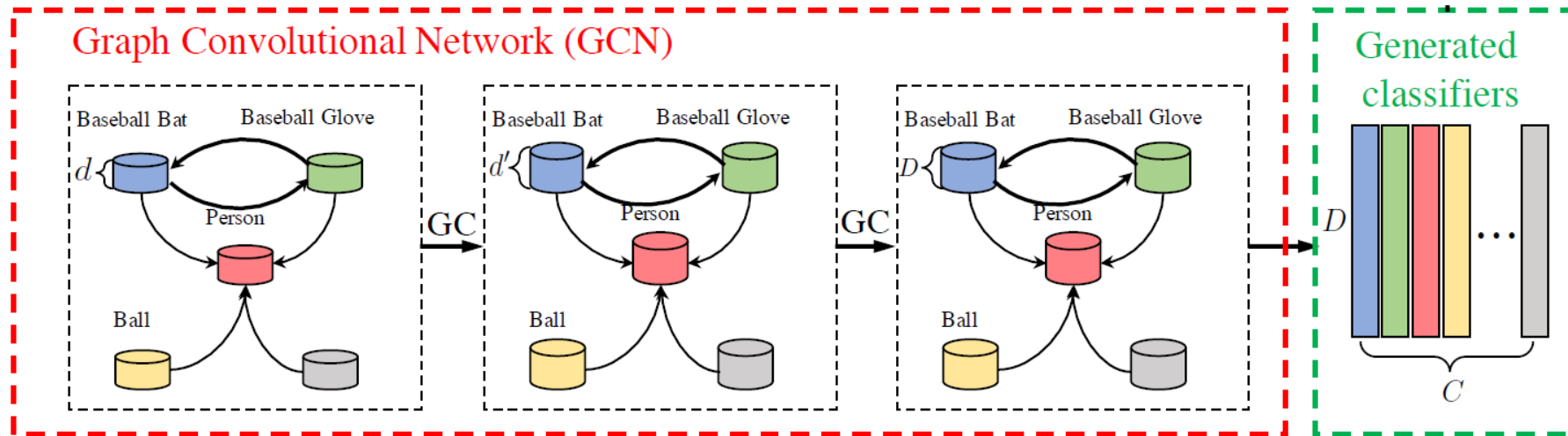
means that when Label<sub>A</sub> appears, Label<sub>B</sub> is likely to appear, but the reverse may not be true



# Method

**Q2:** How to effectively explore these label correlations to improve the classification performance?

Label  $\longrightarrow$  Word Embedding

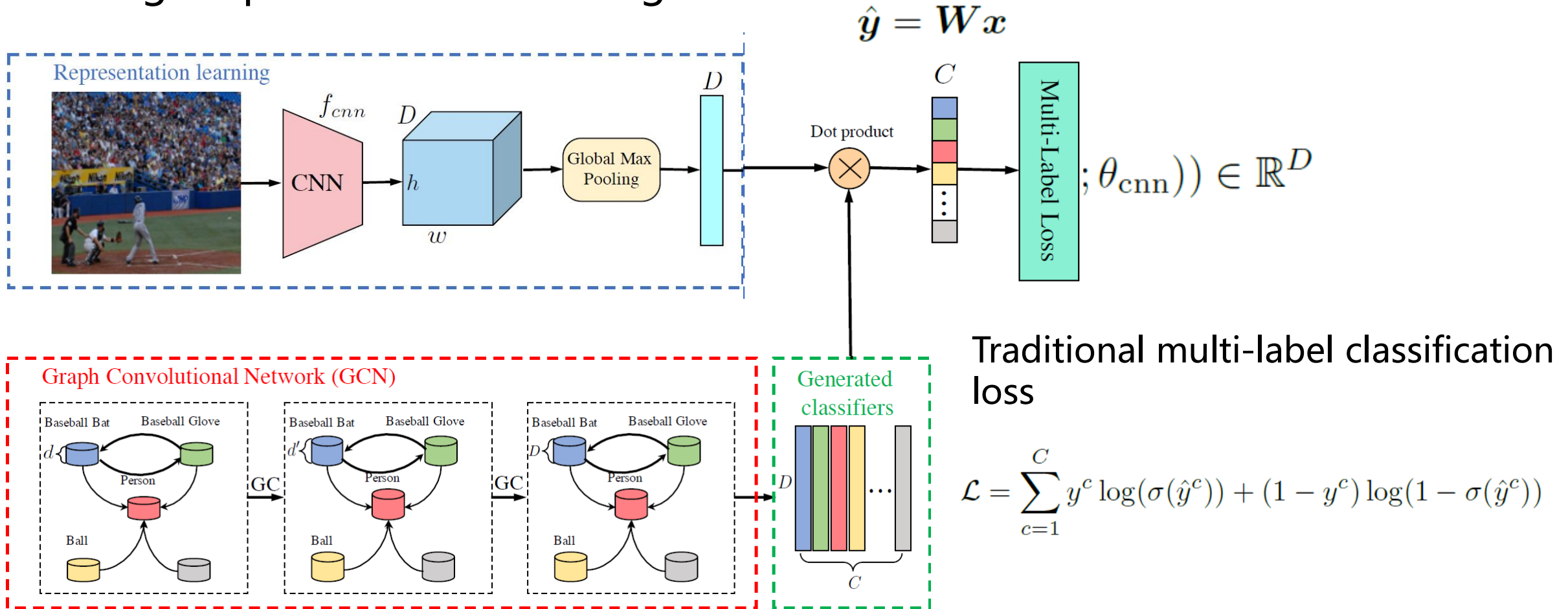


$$H^{l+1} = f(H^l, A) \xrightarrow{\text{Conv}} H^{l+1} = h(\hat{A}H^lW^l)$$

Node  $\downarrow$  Transformation matrix to be learned

# Method

## Image representation learning



# Experiments

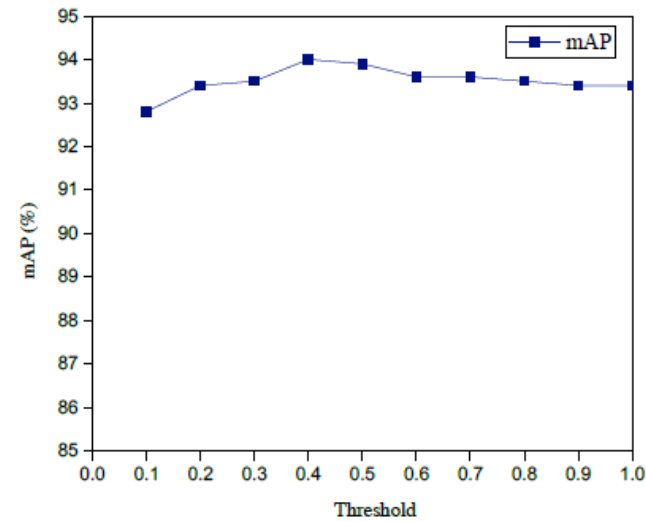
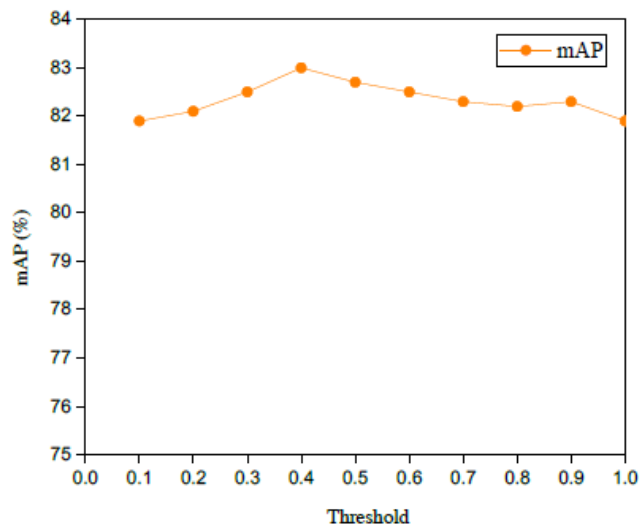
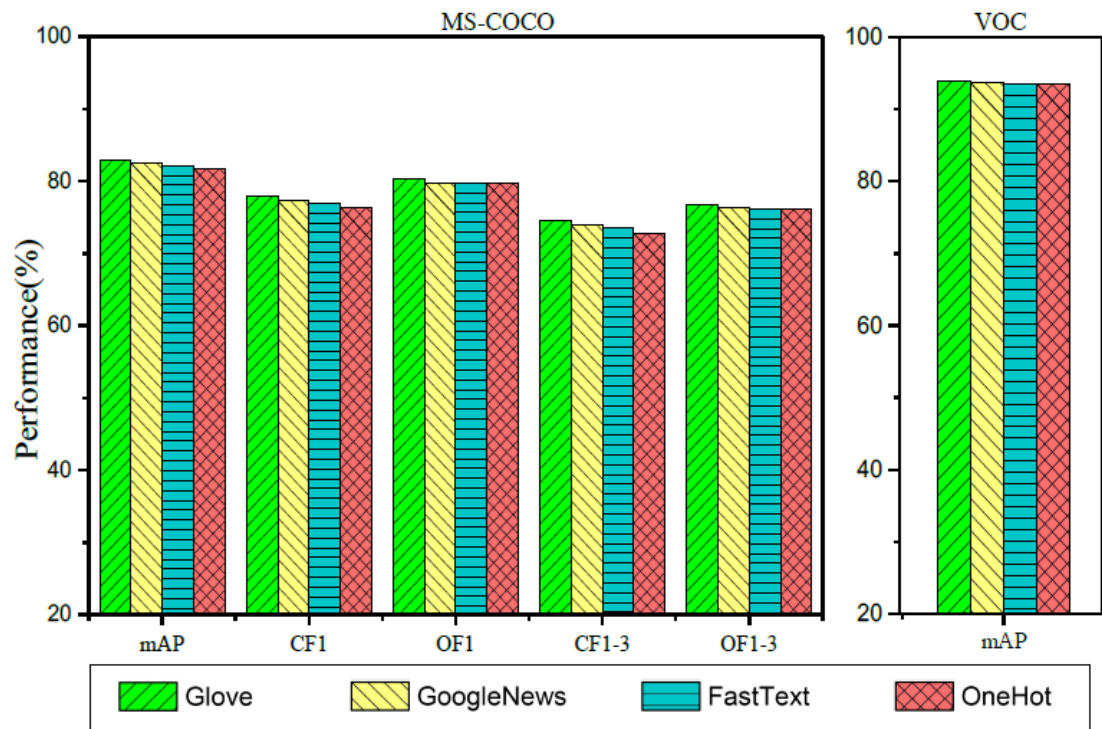
Methods	All							Top-3					
	mAP	CP	CR	CF1	OP	OR	OF1	CP	CR	CF1	OP	OR	OF1
CNN-RNN [28]	61.2	–	–	–	–	–	–	66.0	55.6	60.4	69.2	66.4	67.8
RNN-Attention [29]	–	–	–	–	–	–	–	79.1	58.7	67.4	84.0	63.0	72.0
Order-Free RNN [1]	–	–	–	–	–	–	–	71.6	54.8	62.1	74.2	62.2	67.7
ML-ZSL [15]	–	–	–	–	–	–	–	74.1	<b>64.5</b>	69.0	–	–	–
SRN [36]	77.1	81.6	65.4	71.2	82.7	69.9	75.8	85.2	58.8	67.4	87.4	62.5	72.9
ResNet-101 [10]	77.3	80.2	66.7	72.8	83.9	70.8	76.8	84.1	59.4	69.7	89.1	62.8	73.6
Multi-Evidence [6]	–	80.4	70.2	74.9	85.2	72.5	78.4	84.5	62.2	70.6	89.1	64.3	74.7
ML-GCN (Binary)	80.3	81.1	70.1	75.2	83.8	74.2	78.7	84.9	61.3	71.2	88.8	65.2	75.2
ML-GCN (Re-weighted)	<b>83.0</b>	<b>85.1</b>	<b>72.0</b>	<b>78.0</b>	<b>85.8</b>	<b>75.4</b>	<b>80.3</b>	<b>89.2</b>	64.1	<b>74.6</b>	<b>90.5</b>	<b>66.5</b>	<b>76.7</b>

## MS-COCO

Methods	<i>aero</i>	<i>bike</i>	<i>bird</i>	<i>boat</i>	<i>bottle</i>	<i>bus</i>	<i>car</i>	<i>cat</i>	<i>chair</i>	<i>cow</i>	<i>table</i>	<i>dog</i>	<i>horse</i>	<i>motor</i>	<i>person</i>	<i>plant</i>	<i>sheep</i>	<i>sofa</i>	<i>train</i>	<i>tv</i>	mAP
CNN-RNN [28]	96.7	83.1	94.2	92.8	61.2	82.1	89.1	94.2	64.2	83.6	70.0	92.4	91.7	84.2	93.7	59.8	93.2	75.3	<b>99.7</b>	78.6	84.0
RLSD [34]	96.4	92.7	93.8	94.1	71.2	92.5	94.2	95.7	74.3	90.0	74.2	95.4	96.2	92.1	97.9	66.9	93.5	73.7	97.5	87.6	88.5
VeryDeep [26]	98.9	95.0	96.8	95.4	69.7	90.4	93.5	96.0	74.2	86.6	<b>87.8</b>	96.0	96.3	93.1	97.2	70.0	92.1	80.3	98.1	87.0	89.7
ResNet-101 [10]	99.5	97.7	97.8	96.4	65.7	91.8	96.1	97.6	74.2	80.9	85.0	<b>98.4</b>	96.5	95.9	98.4	70.1	88.3	80.2	98.9	89.2	89.9
FeV+LV [33]	97.9	97.0	96.6	94.6	73.6	93.9	96.5	95.5	73.7	90.3	82.8	95.4	97.7	95.9	98.6	77.6	88.7	78.0	98.3	89.0	90.6
HCP [31]	98.6	97.1	98.0	95.6	75.3	<b>94.7</b>	95.8	97.3	73.1	90.2	80.0	97.3	96.1	94.9	96.3	78.3	94.7	76.2	97.9	91.5	90.9
RNN-Attention [29]	98.6	97.4	96.3	96.2	75.2	92.4	96.5	97.1	76.5	92.0	87.7	96.8	97.5	93.8	98.5	81.6	93.7	82.8	98.6	89.3	91.9
Atten-Reinforce [2]	98.6	97.1	97.1	95.5	75.6	92.8	96.8	97.3	78.3	92.2	87.6	96.9	96.5	93.6	98.5	81.6	93.1	83.2	98.5	89.3	92.0
VGG (Binary)	98.3	97.1	96.1	96.7	75.0	91.4	95.8	95.4	76.7	92.1	85.1	96.7	96.0	95.3	97.8	77.4	93.1	79.7	97.9	89.3	91.1
VGG (Re-weighted)	99.4	97.4	98.0	97.0	77.9	92.4	96.8	97.8	80.8	93.4	87.2	98.0	97.3	95.8	98.8	79.4	95.3	82.2	99.1	91.4	92.8
ML-GCN (Binary)	99.6	98.3	97.9	97.6	78.2	92.3	97.4	97.4	79.2	94.4	86.5	97.4	97.9	97.1	98.7	84.6	95.3	83.0	98.6	90.4	93.1
ML-GCN (Re-weighted)	99.5	<b>98.5</b>	<b>98.6</b>	<b>98.1</b>	80.8	94.6	<b>97.2</b>	98.2	<b>82.3</b>	<b>95.7</b>	86.4	98.2	<b>98.4</b>	<b>96.7</b>	<b>99.0</b>	<b>84.7</b>	<b>96.7</b>	<b>84.3</b>	98.9	<b>93.7</b>	<b>94.0</b>

## VOC 2007

# Experiments

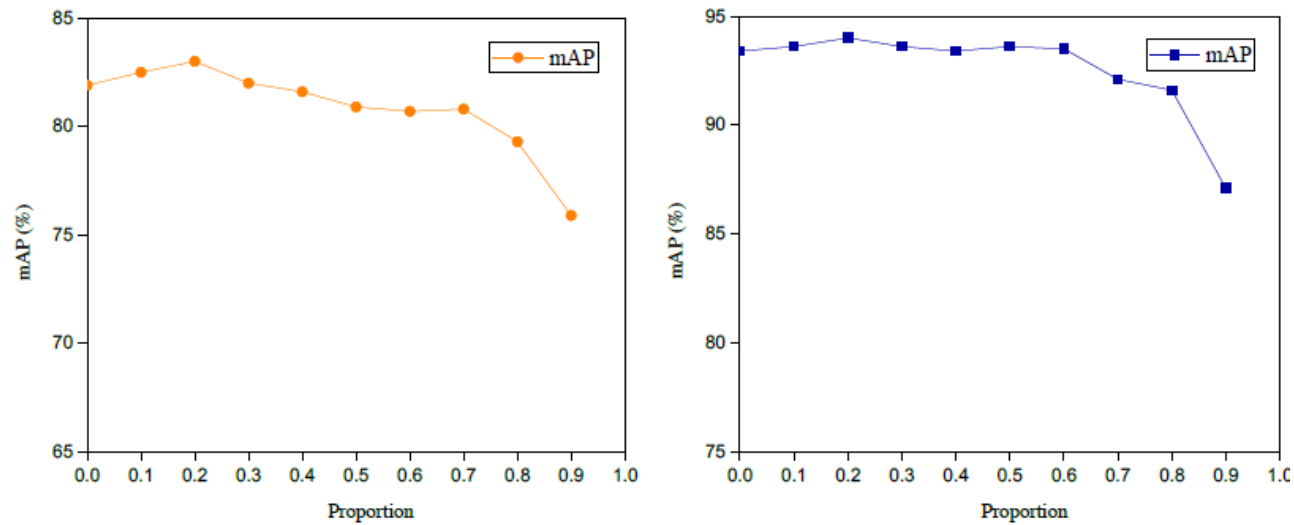


(a) Comparisons on MS-COCO.

(b) Comparisons on VOC 2007.

Figure 5. Accuracy comparisons with different values of  $\tau$ .

# Experiments



(a) Comparisons on MS-COCO.

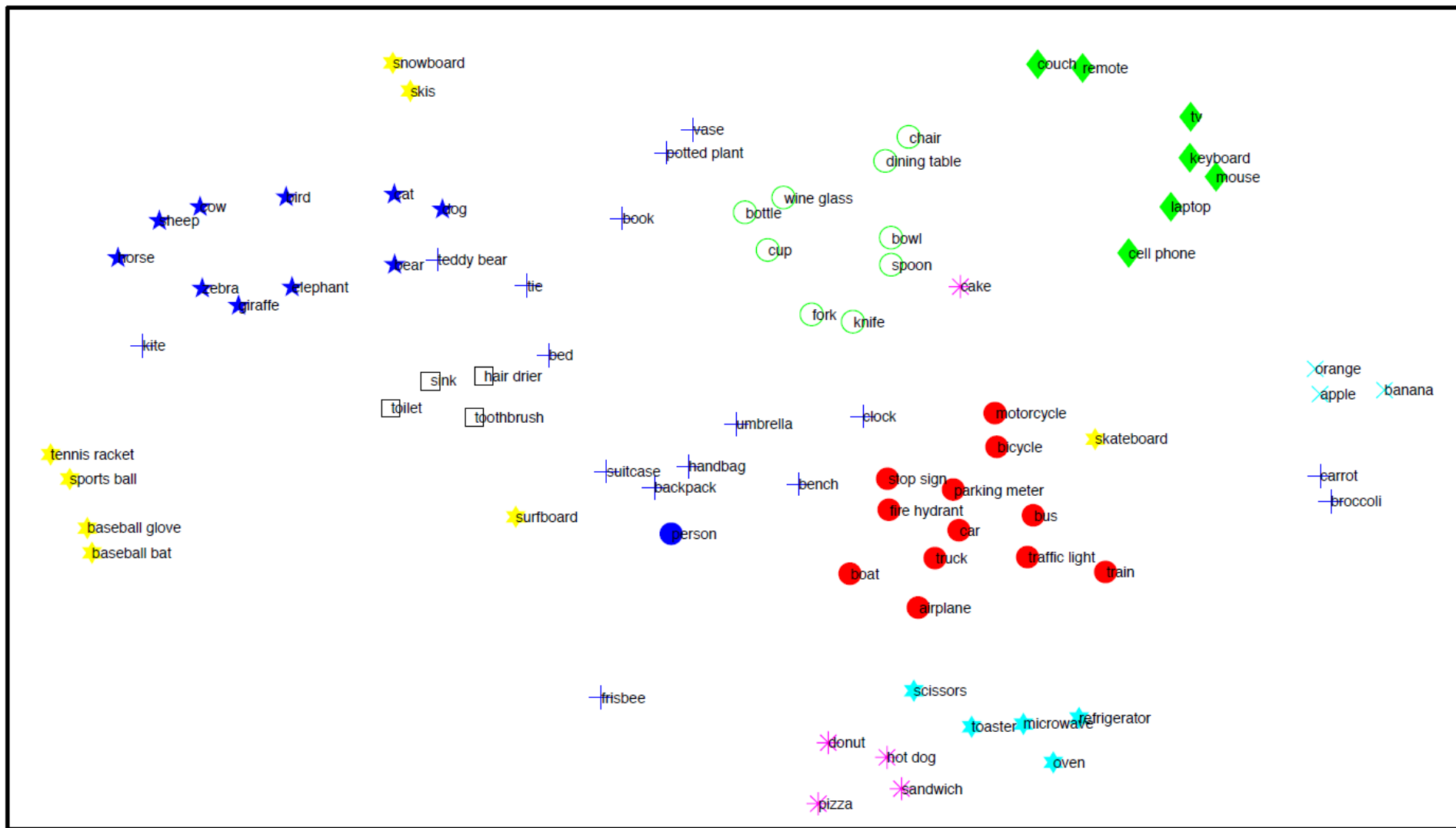
(b) Comparisons on VOC 2007.

Figure 6. Accuracy comparisons with different values of  $p$ . Note that, when  $p = 1$ , the model does not converge.

Table 3. Comparisons with different depths of GCN in our model.

# Layer	MS-COCO					VOC
	All			Top-3		All
	mAP	CF1	OF1	CF1	OF1	mAP
2-layer	<b>83.0</b>	<b>78.0</b>	<b>80.3</b>	<b>74.6</b>	<b>76.7</b>	<b>94.0</b>
3-layer	82.1	76.9	79.7	73.7	76.2	93.6
4-layer	81.1	76.4	79.4	72.5	75.8	93.0

# Experiments



(a) t-SNE on the learned inter-dependent classifiers by our model.

# Experiments

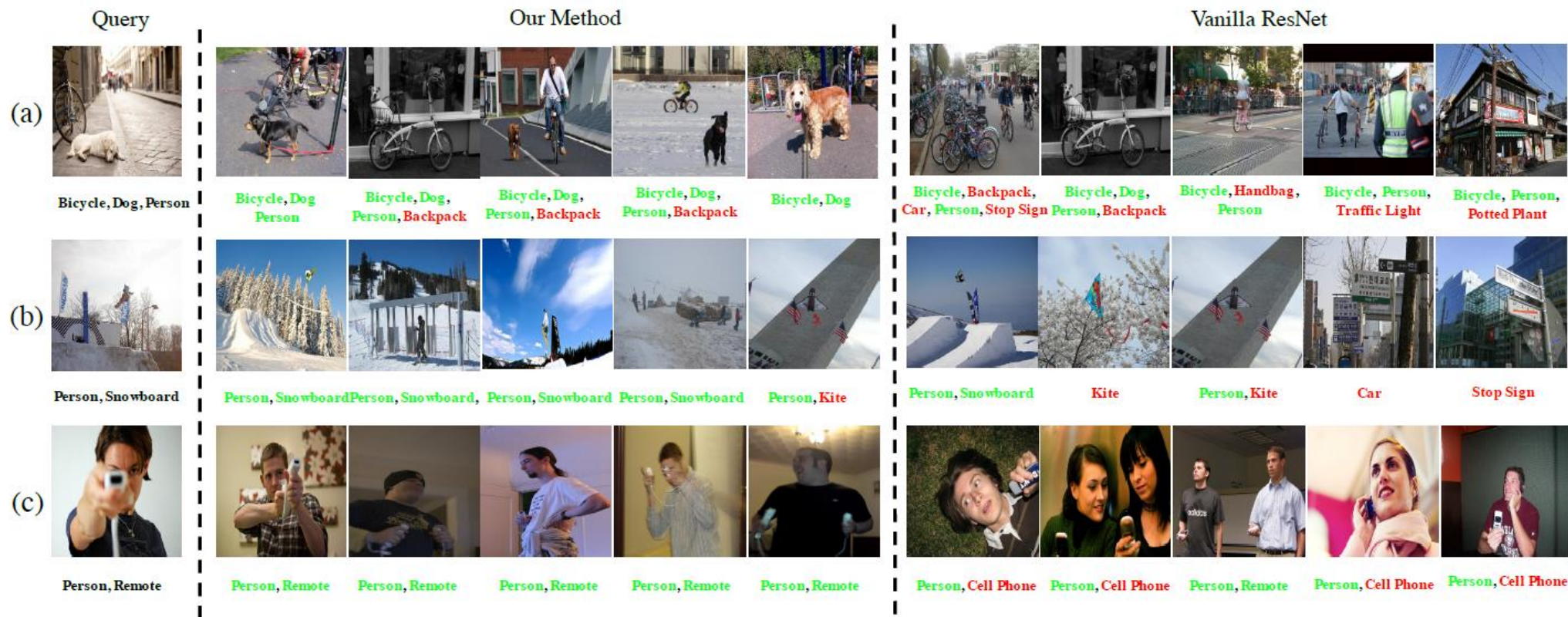


Figure 7. Top-5 returned images with the query image. The returned results on the left are based on our proposed ML-GCN, while the results on the right are vanilla ResNet. All results are sorted in the ascending order according to the distance from the query image.

Thanks

---