

# Multi-label Learning via Codewords

2018 ICTAI, Volos, Greece, November 5–7, 2018.

Mahlagha Sedghi<sup>1</sup>   Yinjie Huang<sup>1</sup>   Michael Georgiopoulos<sup>1</sup>  
Georgios C. Anagnostopoulos<sup>2</sup>

<sup>1</sup>Machine Learning Laboratory, University of Central Florida, US

<sup>2</sup>ICE Laboratory, Florida Institute of Technology, US

November 05<sup>th</sup>, 2018

# Table of Contents

- 1 Introduction
- 2 Formulation
- 3 Algorithm
- 4 Experiments
- 5 Summary
- 6 References

# Section 1

## Introduction

# Classification Paradigms

Pick one

Label 1	✓
Label 2	

**Binary**

Pick one

Label 1	
Label 2	
Label 3	
Label 4	✓
...	
...	
Label L	

**Multi-class**

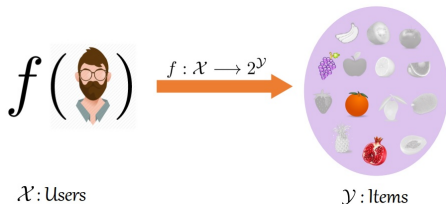
Pick all applicable

Label 1	
Label 2	✓
Label 3	
Label 4	✓
...	
...	
Label L	✓

**Multi-label**

# Multi-label Classification Examples

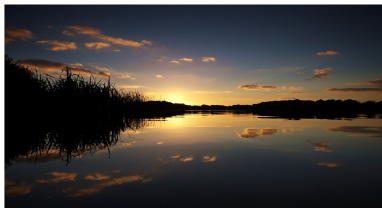
What items would this user buy?



Who are present in this selfie?



What tags are related to this image?



# Label Correlations

- In most realistic settings, presence or absence of the labels are correlated.
  - In image annotation, if an image features a *fish* tag, it will most likely also feature an *ocean* tag, rather than a *sky* tag.
- Structured Output Prediction (SOP) approaches can be used to capture these correlations.
  - Exact SOP training algorithms applicable to Multi-label problems suffer from computationally intensive training time. ( $\mathcal{O}(2^L)$ , with  $L$  as number of labels)

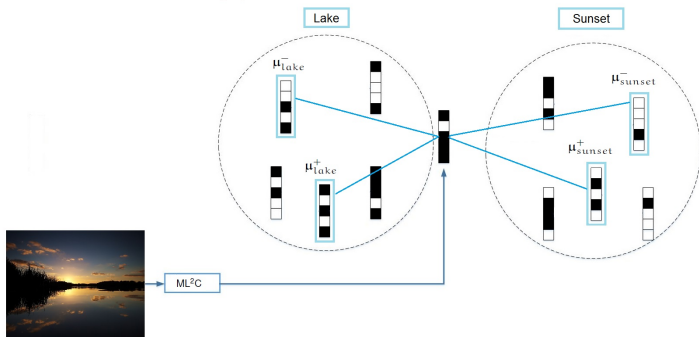
## Section 2

# Formulation

# Formulation

Multi-label Learning via Codewords (ML<sup>2</sup>C) utilizes

- Hash function/encoders to map the input data into a common  $B$ -bit code space.
- $L$  pairs of codewords in the same space, where each pair is used for presence/absence of a particular label.
- The prediction of each label is based on the Hamming distance of the input's hash code with the codeword pair.



# Formulation

- Data:  $\{(\mathbf{x}_n, \mathbf{y}_n)\}_{n \in \mathbb{N}_N}$ , where  $\mathbf{x}_n \in \mathcal{X} \triangleq \mathbb{R}^D$  and  $\mathbf{y}_n \in \mathbb{H} \triangleq \{-1, 1\}^L$
- SOPs are performed in the hash space.
  - True labels are represented by codewords set  $\mathbf{M} \triangleq [\boldsymbol{\mu}_1(1), \boldsymbol{\mu}_1(-1), \dots, \boldsymbol{\mu}_L(1), \boldsymbol{\mu}_L(-1)] \in \mathbb{R}^{B \times 2L}$
  - Weight parameters:  $\mathbf{W} \triangleq [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_B] \in \mathbb{R}^{D \times B}$ , where  $\mathbf{w}_b \in \mathbb{R}^D$  corresponds to each bit.

- Hash code: 
$$\mathbf{h}(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathbb{H}^L} f(\mathbf{x}, \mathbf{y}) \quad (1)$$

- Scoring function for each label  $l \in \mathbb{N}_L$ :

$$f_b(\mathbf{x}_n, h_b; \mathbf{w}_b) = \langle \mathbf{w}_b, \phi_b(\mathbf{x}_n, h_b) \rangle \quad \forall b \in \mathbb{N}_B$$

$$f(\mathbf{x}_n, \mathbf{h}; \mathbf{W}) = \mathbf{x}_n^T \mathbf{W} \mathbf{h} \quad (2)$$

- Feature map for each label  $l \in \mathbb{N}_L$ :

$$\phi_b(\mathbf{x}_n, h_b) = h_b \mathbf{x}_n$$

$$\phi(\mathbf{x}_n, \mathbf{h}) = \mathbf{x}_n \mathbf{h}^T \quad (3)$$

# Formulation

- By adjusting parameters  $W$  and  $M$ ,  $ML^2C$  attempts to reduce the distortion measure:

$$E(\mathbf{W}, \mathbf{M}) \triangleq \sum_{n \in \mathbb{N}_N} \sum_{l \in \mathbb{N}_L} [y_{nl} = 1] d_H(\mu_l(1), \mathbf{h}(\mathbf{x}_n)) \\ + [y_{nl} = -1] d_H(\mu_l(-1), \mathbf{h}(\mathbf{x}_n)) \quad (4)$$

where  $d_H$  is the Hamming distance defined as

$$d_H(\mu_l(y_{nl}), \mathbf{h}(\mathbf{x}_n)) \triangleq \sum_b [\mu_{lb}(y_{nl}) h_b(\mathbf{x}_n) < 0], b \in \mathbb{N}_B \quad (5)$$

- Label prediction in the original space

$$\hat{y}_l(\mathbf{x}_n) = \arg \min_{y_l \in \{-1, 1\}} d_H(\mu_l(y_l), \mathbf{h}(\mathbf{x}_n)), \forall l \in \mathbb{N}_L \quad (6)$$

## Section 3

# Algorithm

# Algorithm - Observations

- We decompose the error of the final predicted label  $\hat{\mathbf{y}}$  to that of each label  $l$ :  $Err(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{l \in \mathbb{N}_L} err(y_l, \hat{y}_l)$
- Based on 4, the label-wise error can be defined as

$$err(y_l, \hat{y}_l) \equiv \Delta_l(\boldsymbol{\mu}_l(y_l), \hat{\mathbf{h}}) = [y_l = 1]d_H(\boldsymbol{\mu}_l(1), \hat{\mathbf{h}}) + [y_l = -1]d_H(\boldsymbol{\mu}_l(-1), \hat{\mathbf{h}}) \quad (7)$$

- This problem is non-convex, and hence difficult to optimize directly.
- Upperbounding the error function with a convex loss over parameters could be helpful.
- A two-step Block Coordinate Descent algorithm would optimize the surrogate loss.

# Algorithm - Surrogate Loss

- The error defined in Eq. (7) has the following surrogate loss:

$$\Delta_l(\mu_l(y_l), \hat{\mathbf{h}}) \leq \sup_{\mathbf{h}' \in \mathbb{H}^B} \Delta_l(\mu_l(y_l), \mathbf{h}') + f(\mathbf{x}_n, \mathbf{h}'; \mathbf{W}) - f(\mathbf{x}_n, \mu_l(y_l); \mathbf{W}) \quad (8)$$

- Considering a regularization term on  $\mathbf{W}$ , one can alternatively solve the following problem:

$$\min_{\mathbf{W}, \mathbf{M}} \sum_{n \in \mathbb{N}_{\mathcal{N}}} \sum_{l \in \mathbb{N}_L} \sup_{\mathbf{h}' \in \mathbb{H}^B} \{ \Delta_l(\mu_l(y_l), \mathbf{h}') + f(\mathbf{x}_n, \mathbf{h}'; \mathbf{W}) - f(\mathbf{x}_n, \mu_l(y_l); \mathbf{W}) + \frac{\lambda}{2} \|\mathbf{W}\|_F^2 \} \quad (9)$$

# Algorithm - First block minimization

- Prob. (10) is equivalent to Prob. (9) through
  - Defining slack variables  $\xi_{nl}$  which represent maximal value of the surrogate loss,
  - Reforming the surrogate loss as constraints

$$\min_{\mathbf{W}} \sum_{n \in \mathbb{N}_{\mathcal{N}}} \sum_{l \in \mathbb{N}_L} \xi_{nl} + \frac{\lambda}{2} \|\mathbf{W}\|_F^2 \quad (10)$$

$$\text{s.t. } f(\mathbf{x}_n, \mu_l(y_l); \mathbf{W}) - f(\mathbf{x}_n, \mathbf{h}'; \mathbf{W}) \geq \Delta_l(\mu_l(y_l), \mathbf{h}') - \xi_{nl} \quad \forall \mathbf{h}' \in \mathbb{H}^B$$

- Prob. (10) is easily recognized as the standard formulation of Structured Support Vector Machine (SSVM) ( $\mathcal{O}(L2^B)$ ).

# Algorithm - Second block minimization

- Codewords Optimization Step:

$$\min_{\mathbf{M}} \sum_{n \in \mathbb{N}_{\mathcal{N}}} \sum_{l \in \mathbb{N}_L} \sup_{\mathbf{h}' \in \mathbb{H}^B} \{ \Delta_l(\mu_l(y_l), \mathbf{h}') + f(\mathbf{x}_n, \mathbf{h}'; \mathbf{W}) - f(\mathbf{x}_n, \mu_l(y_l); \mathbf{W}) \} \quad (11)$$

- Having  $W$  fixed, the supremum value of the surrogate loss w.r.t.  $\mathbf{h}'$  is attained by  $\tilde{\mathbf{h}}$ , calculated in the previous block.
- The problem is decomposable to  $L$  separate optimization problems in terms of  $\mu_l(y_l)$ , which can be solved in parallel:

$$\min_{\mu_l(y_l)} \sum_{n \in \mathbb{N}_{\mathcal{N}}} [y_{nl} = 1][\mu_l(1)\tilde{\mathbf{h}}_n < 0] + [y_{nl} = -1][\mu_l(-1)\tilde{\mathbf{h}}_n < 0] \\ + \mathbf{w}_b^T \mathbf{x}_n (\tilde{\mathbf{h}}_n - [y_{nl} = 1]\mu_l(1) - [y_{nl} = -1]\mu_l(-1)) \quad (12)$$

- Prob. (12) can be further simplified in terms of bit-wise quantities, and is solved by simple substitution ( $\mathcal{O}(LB)$ ).

# Algorithm - Optimization of Prob. (4)

**Input:** Codeword Bit Length  $B$ , Training Samples  $X$  with multi-labels  $Y$ .

**Output:**  $W$ ,  $M$ .

01. Initialize  $W$ ,  $M$ .

02. **While Not Converged**

03. Weight Optimization Step:

04. **For each label**

05. Score Function  $\leftarrow$  Eq. (2).

06. Feature Map  $\leftarrow$  Eq. (3).

07. Cost Function  $\leftarrow$  Eq. (7).

08. Solve SSVM with `SSVM-MATLAB`.

09. **End For**

10. Codewords Optimization Step: Optimize over Prob. (11).

11. **End While**

The algorithm has  $\mathcal{O}(L2^B)$  overall complexity, which is computationally superior to the state-of-the-art.

## Section 4

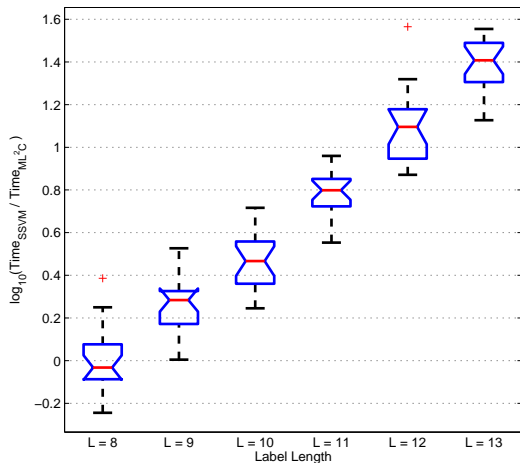
# Experiments

# Datasets

A subset of following datasets are used for different experiments.

Dataset	#Inst	#Dim	#L	Card	Den	Dist	Dom
<b>Birds</b>	645	260	19	1.014	0.053	133	Audio
<b>CAL500</b>	502	68	174	26.044	0.150	502	Music
<b>Emotions</b>	593	72	6	1.869	0.311	27	Music
<b>Scene</b>	2407	294	6	1.074	0.179	15	Image
<b>Yeast</b>	2417	103	14	4.237	0.303	198	Biology
<b>Rcv1v2-s1</b>	6000	47236	101	2.880	0.029	1028	Text
<b>Rcv1v2-s2</b>	6000	47236	101	2.634	0.026	954	Text
<b>Rcv1v2-s3</b>	6000	47236	101	2.614	0.026	939	Text
<b>Rcv1v2-s4</b>	6000	47236	101	2.484	0.025	816	Text
<b>Rcv1v2-s5</b>	6000	47236	101	2.642	0.026	946	Text

# Training Time Comparison



**Figure:** Comparison of training time for ML<sup>2</sup>C and SSVM on the *Birds* data set. ML<sup>2</sup>C used a code word length of  $B = 5$ .

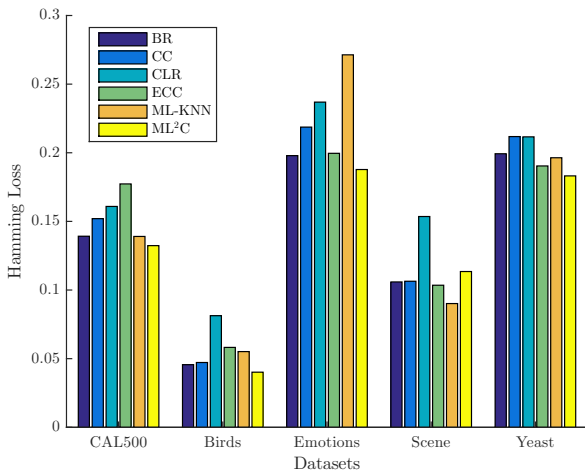
# Performance Comparison - Methods and Measures

- Methods to compare with:
  - Classifier Chains (CC) [Read et al., 2011]
  - Calibrated Label Ranking (CLR) [Fürnkranz et al., 2008]
  - Ensemble of Classifier Chains (ECC) [Read et al., 2011]
  - Multi-Label K Nearest Neighbor (ML-KNN) [Zhang and Zhou, 2007]
  - Binary Relevance - Support Vector Machine (BR-SVM) [Godbole and Sarawagi, 2004]

- Evaluation Measures:

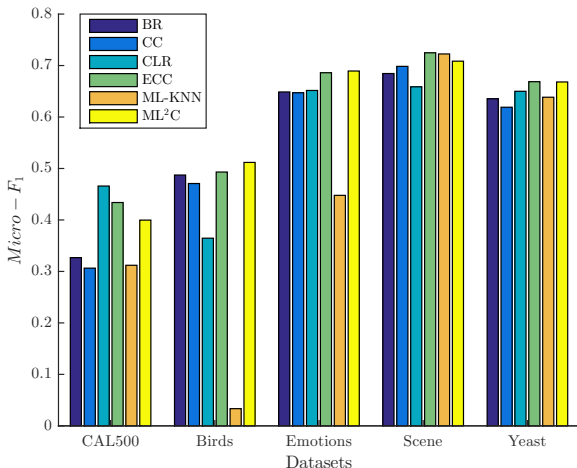
- Hamming Loss  $\triangleq \frac{1}{N} \sum_{n \in \mathbb{N}_N} |\mathbf{y}_n \cap \hat{\mathbf{y}}_n|$
- *micro* -  $F_1 \triangleq \frac{\sum_{l \in \mathbb{N}_L} 2TP_l}{\sum_{l \in \mathbb{N}_L} 2TP_l + FP_l + FN_l}$
- *macro* -  $F_1 \triangleq \frac{1}{L} \sum_{l \in \mathbb{N}_L} \frac{2TP_l}{2TP_l + FP_l + FN_l}$

# Performance Comparison - Hamming Loss



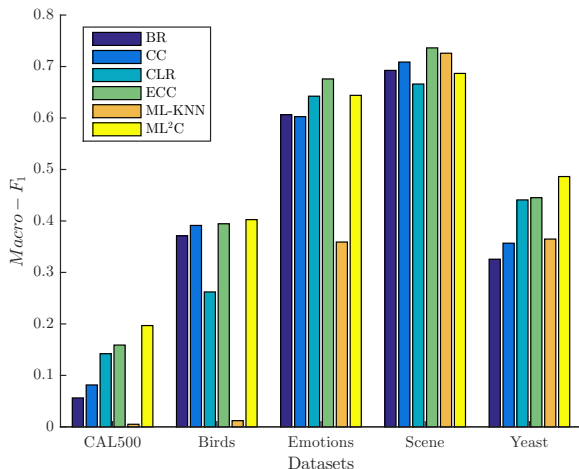
**Figure:** Average Hamming Loss of 6 algorithms over 5 medium-size datasets. The smaller the Hamming Loss is, the better the algorithms perform.

# Performance Comparison - *micro* - $F_1$



**Figure:** Average *micro* -  $F_1$  of 6 algorithms over 5 medium-size datasets. The larger the  $F_1$  measures are, the better the algorithms perform.

# Performance Comparison - *macro* – $F_1$



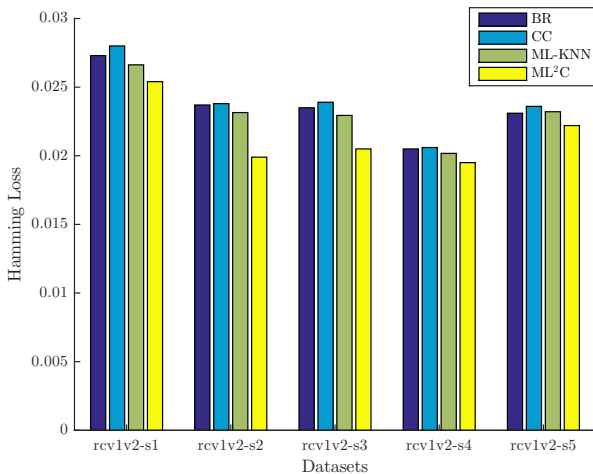
**Figure:** Average *macro* –  $F_1$  of 6 algorithms over 5 medium-size datasets. The larger the  $F_1$  measures are, the better the algorithms perform.

# Performance Comparison - Ranking

**Table:** Ranking results of all algorithms on 5 benchmark datasets of medium size. For each dataset, the statistically best and comparable results for a family-wise error significance level of 0.05 are highlighted in boldface. All algorithms are ranked from best to worst. If the performance is statistically comparable, the algorithms share the same rank.

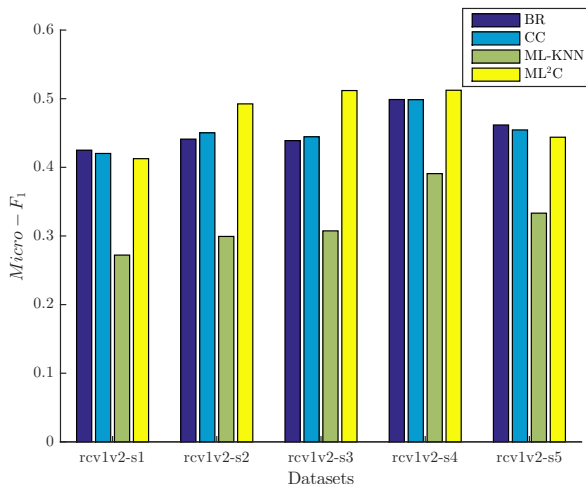
Dataset	ML <sup>2</sup> C	BR-SVM	ML-KNN	ECC	CLR	CC
<b>CAL500</b>	<b>1st</b>	<i>2nd</i>	<i>2nd</i>	<i>5th</i>	<i>4th</i>	<i>3rd</i>
<b>Birds</b>	<b>1st</b>	<i>2nd</i>	<i>2nd</i>	<b>1st</b>	<i>3rd</i>	<b>1st</b>
<b>Emotions</b>	<b>1st</b>	<i>2nd</i>	<i>5th</i>	<i>2nd</i>	<i>4th</i>	<i>3rd</i>
<b>Scene</b>	<i>3rd</i>	<b>1st</b>	<i>2nd</i>	<i>2nd</i>	<b>1st</b>	<i>2nd</i>
<b>Yeast</b>	<b>1st</b>	<i>3rd</i>	<i>2nd</i>	<b>1st</b>	<i>4th</i>	<i>4th</i>

# Performance Comparison - Hamming Loss



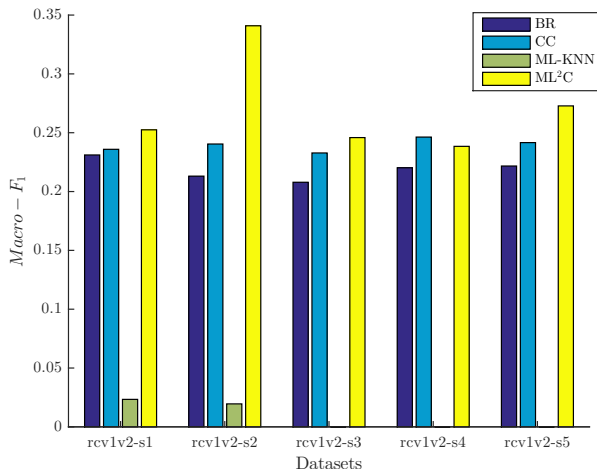
**Figure:** Average Hamming Loss of 6 algorithms over 5 large-size datasets. The smaller the Hamming Loss is, the better the algorithms perform.

# Performance Comparison - *micro* - $F_1$



**Figure:** Average *micro* -  $F_1$  of 6 algorithms over 5 large-size datasets. The larger the  $F_1$  measures are, the better the algorithms perform.

# Performance Comparison - *macro* – $F_1$



**Figure:** Average *macro* –  $F_1$  of 6 algorithms over 5 large-size datasets. The larger the  $F_1$  measures are, the better the algorithms perform.

# Performance Comparison - Ranking

**Table:** Ranking results of algorithms on 5 benchmark datasets of large size. For each dataset, the statistically best and comparable results for a family-wise error significance level of 0.05 are highlighted in boldface. All algorithms are ranked from best to worst. If the performance is statistically comparable, the algorithms share the same rank.

Data Sets	ML <sup>2</sup> C	BR-SVM	ML-KNN	CC
Rcv1v2-s1	<b>1st</b>	<i>3rd</i>	<i>2nd</i>	<i>3rd</i>
Rcv1v2-s2	<b>1st</b>	<i>3rd</i>	<i>2nd</i>	<i>3rd</i>
Rcv1v2-s3	<b>1st</b>	<i>3rd</i>	<i>2nd</i>	<i>3rd</i>
Rcv1v2-s4	<b>1st</b>	<i>2nd</i>	<b>1st</b>	<i>2nd</i>
Rcv1v2-s5	<b>1st</b>	<i>2nd</i>	<i>2nd</i>	<i>3rd</i>

# Image Annotation

**Table:** Street scene of the SUN dataset. Wrong annotations (red), missed annotations (inside parentheses).



<b>Ground Truth</b>	Car, Building, Sidewalk, Road, Sky, Tree	Car, Building, Sidewalk, Road, Sky, Tree	Car, Building, Sidewalk, Road, Person, Sky, Tree
<b>ML<sup>2</sup>C</b>	Car, Building, Sidewalk, Road, Sky, Tree	Car, Building, Sidewalk, Road, Sky, Tree	Car, Building, Sidewalk, Road, Person, Sky, Tree
<b>BR-SVM</b>	Car, Building, Sidewalk, Road, <b>Person</b> , Sky, ( <b>Tree</b> )	Car, Building, Sidewalk, Road, <b>Person</b> , Sky, ( <b>Tree</b> )	Car, Building, Sidewalk, Road, Person, Sky, ( <b>Tree</b> )
<b>ML-KNN</b>	Car, Building, Sidewalk, Road, <b>Person</b> , Sky, ( <b>Tree</b> )	Car, Building, Sidewalk, Road, <b>Person</b> , Sky, ( <b>Tree</b> )	Car, Building, Sidewalk, Road, Person, Sky, ( <b>Tree</b> )
<b>ECC</b>	Car, Building, ( <b>Tree</b> ), <b>Window</b> , Sidewalk, Road, Sky	Car, Building, Sidewalk, Road, Sky, ( <b>Tree</b> )	Car, Building, Sidewalk, Road, ( <b>Person</b> ), Sky, Tree
<b>CLR</b>	Car, Building, Tree, <b>Window</b> , Sidewalk, Road, Sky	Car, Building, Sidewalk, Road, Sky, ( <b>Tree</b> )	Car, Building, Sidewalk, Road, ( <b>Person</b> ), Sky, Tree
<b>CC</b>	Car, Building, Sidewalk, Road, <b>Window</b> , Sky, ( <b>Tree</b> )	Car, Building, Sidewalk, Road, <b>Person</b> , Sky, ( <b>Tree</b> )	Car, Building, Sidewalk, Road, ( <b>Person</b> ), Sky, Tree

# Image Annotation

**Table:** *Mountain* scene of the *SUN* dataset. Wrong annotations (red), missed annotations (inside parentheses).



<b>Ground Truth</b>	Mountain, Sky	Mountain, Sky, Snow	Mountain, Sky
<b>ML<sup>2</sup>C</b>	Mountain, Sky	Mountain, Sky, (Snow)	Mountain, Sky
<b>BR-SVM</b>	Mountain, Sky, Snow	Mountain, Sky, (Tree), Snow	Mountain, Sky, Tree, Snow
<b>ML-KNN</b>	Mountain, Sky, Tree	Mountain, Sky, (Tree), (Snow)	Mountain, Sky, Tree
<b>ECC</b>	Mountain, Sky, Tree, Snow	Mountain, Sky, Tree, (Snow)	Mountain, Sky, Tree
<b>CLR</b>	Mountain, Sky	Mountain, Sky, (Snow)	Mountain, Sky, Tree
<b>CC</b>	Mountain, Sky, Snow	Mountain, Sky, Tree, (Snow)	Mountain, Sky, Tree

# Section 5

## Summary

# Summary

- A novel multi-label learning framework - Multi-label Learning via Codewords ( $ML^2C$ ) is proposed.
- $ML^2C$  reduces the training time complexity from typical Structured Output Prediction (SOP) approaches ( $\mathcal{O}(2^L)$ ) to  $\mathcal{O}(L2^B)$ .
- Its training algorithm is simple to implement.
- Experiments on 10 benchmark datasets, compared with 5 other state-of-art methods, show  $ML^2C$  is superior.

# Thank You!

Thanks for your time.

## Section 6

### References

# References I

-  Fürnkranz, J., Hüllermeier, E., LozaMencía, E., and Brinker, K. (2008). Multilabel classification via calibrated label ranking. *Machine Learning*, 73(2):133–153.
-  Godbole, S. and Sarawagi, S. (2004). Discriminative methods for multi-labeled classification. In *Pacific-Asia conference on knowledge discovery and data mining*, pages 22–30. Springer.
-  Read, J., Pfahringer, B., Holmes, G., and Frank, E. (2011). Classifier chains for multi-label classification. *Machine Learning*, 85:333–359.
-  Zhang, M.-l. and Zhou, Z.-h. (2007). MI-knn: A lazy learning approach to multi-label learning. *Pattern Recognition*, 40:2007.