

# Common and Discriminative Semantic Pursuit for Multi-Modal Multi-Label Learning

Yi Zhang and Jundong Shen and Zhecheng Zhang and Chongjun Wang<sup>1</sup>

**Abstract.** Multi-modal multi-label (MMML) learning provides an important framework to learn complex objects with diverse representations and annotations. Most existing multi-modal multi-label learning approaches focus on exploiting shared information of all modalities, but neglect specific information of each modality. Besides, how to effectively utilize relationship among modalities is also a challenging issue. In this paper, we propose a novel MMML learning approach called Common and Discriminative Semantic Pursuit (CoDiSP), which learns low-dimensional common representation with all modalities, and extracts discriminative information of each modality by enforcing orthogonal constraint. Meanwhile, the common representation is used as a new modality and added to the specific modal sequence. Furthermore, CoDiSP learns deep models with adaptive depth and exploits label correlations simultaneously based on the extracted modal sequence. Finally, extensive experiments on several benchmark MMML datasets show superior performance of CoDiSP compared with other state-of-the-art approaches.

## 1 INTRODUCTION

On the one hand, multi-modal learning has attracted much attention with the rapid development of data collection technology, where data collected from diverse sources can be represented with multiple modal features. In contrast to single modal learning, multi-modal learning mainly exploits the consistent and complementary properties among different modalities and improves the learning performance [26]. The major challenge of multi-modal problem is how to jointly model heterogeneous modalities in a mutually beneficial way.

On the other hand, multi-label learning [33] has been an important and practical research topic, since one instance may have multiple annotations. As a fundamental framework handling objects with multiple annotations, multi-label learning has been widely applied in many real-world applications, such as image annotation [24], document categorization [14], information retrieval [6] and bioinformatics [31]. A challenging issue for multi-label learning is the exploitation of label correlations [5].

In real-world applications, data are often represented with multiple modalities and associated with multiple labels simultaneously. Multi-modal multi-label (MMML) [29] learning provides a learning framework for modeling such complex objects. Existing MMML learning approaches mainly focus on exploiting shared information among all modalities to eliminate noise and redundancy, which map each modality into a shared subspace. However, it is usually done in an independent way, which neglects communication among dif-

ferent modalities. To overcome the above challenges, we present a novel end-to-end multi-modal multi-label neural network framework named Common and Discriminative Semantic Pursuit (CoDiSP), to exploit relationship among different modalities and label correlations with the help of extracted common and specific modal features. The main contributions of this paper are summarized as follows:

- A novel MMML learning approach called Common and Discriminative Semantic Pursuit (CoDiSP) is proposed, which can exploit common and discriminative information, and make full use of inter-relationship among different modalities.
- Without pre-determining the modal order for prediction, our approach is able to sequentially learn the modal dependency and label correlations with the extracted modal sequence.
- Experiments on benchmark datasets demonstrate that CoDiSP performs favorably against state-of-the-art approaches. Besides, we conduct extensive experiments to analyze common semantic, modal dependency and convergence of CoDiSP.

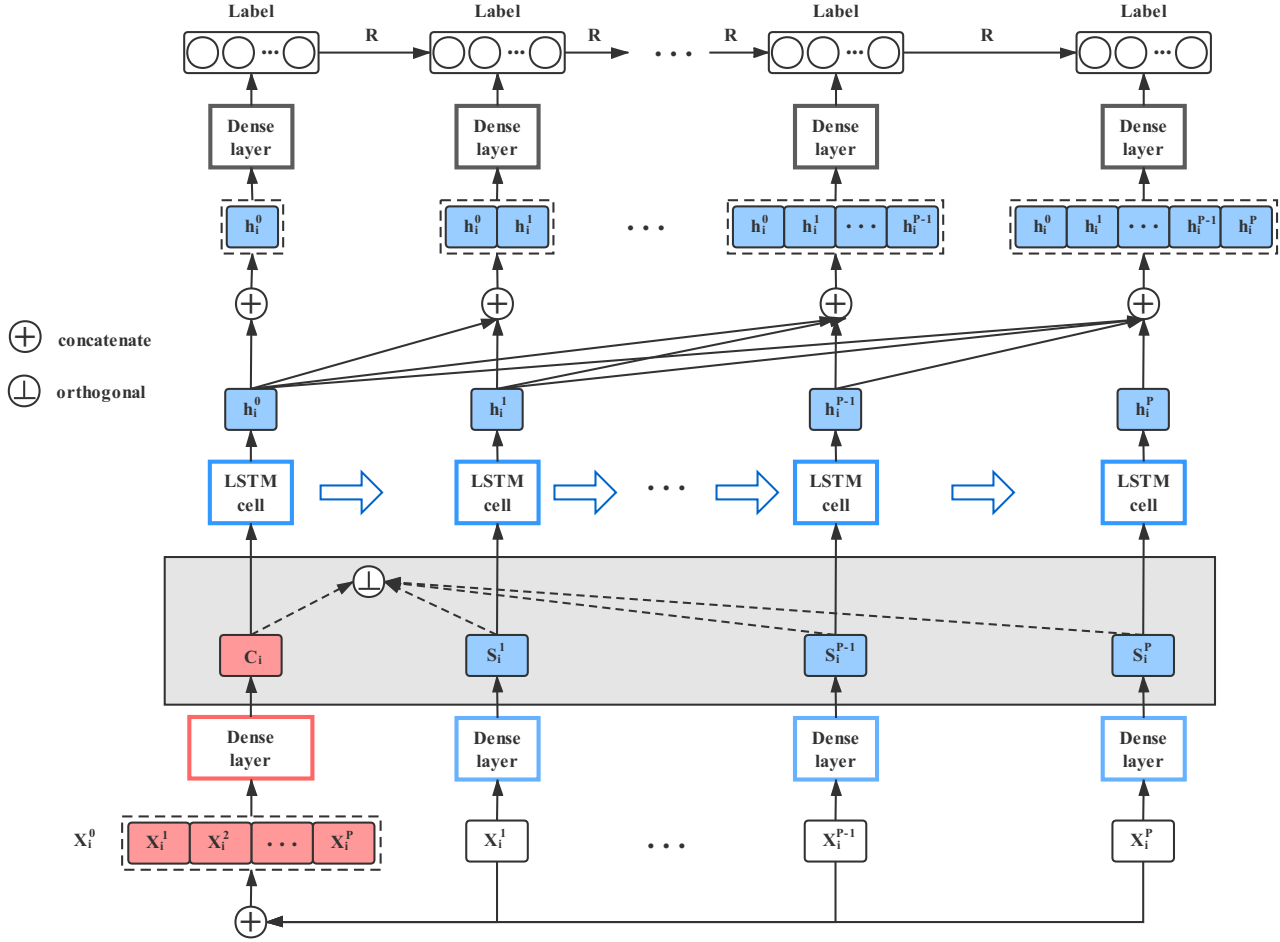
The remainder of this paper is organized as follows. Section 2 briefly reviews some related work of multi-modal multi-label learning. Section 3 presents technical details of the proposed approach. Section 4 reports detailed results of comparative experiments. Finally, Section 5 gives the conclusion.

## 2 RELATED WORK

In this section, we will briefly review some state-of-the-art approaches in both multi-modal and multi-label learning.

For multi-label learning, each instance is associated with multiple interdependent labels and the goal is to exploit various types of label correlations. In terms of the order of label correlations, these approaches can be divided into three strategies. For first-order approach, Binary Relevance (BR) [1] takes each label independently, which neglects the relationship among labels. For second-order approach, Calibrated Label Ranking (CLR) [4] considers the pairwise relationships between labels. For high-order approach, Classifier Chains (CC) [16] addresses connections among random subsets of labels. So far, many approaches have been developed to improve the performance of multi-label learning by exploring various types of label correlations [11] [35]. For example, a boosting approach [10] exploits label correlations with a hypothesis reuse mechanism. However most of the existing approaches take label correlations as prior knowledge [9], which may not correctly characterize the real relationships among labels and final predictions are not explicitly correlated. To tackle this problem, CAMEL [3] learns the high-order label correlations via sparse reconstruction in the label space.

<sup>1</sup> National Key Laboratory for Novel Software Technology at Nanjing University, Nanjing University, Nanjing 210023, China, email: {njuzhangy, jd-shen, zzc}@mail.nju.edu.cn, chjwang@nju.edu.cn



**Figure 1.** The overall flowchart of the proposed CoDiSP approach. Firstly, for  $i$ -th instance  $\mathbf{X}_i = [\mathbf{X}_i^1, \mathbf{X}_i^2, \dots, \mathbf{X}_i^P]$  with  $P$  modalities, we concatenate all modalities as raw common modality  $\mathbf{X}_i^0$ . Secondly, CoDiSP enforces orthogonal constraint to exploit common information  $\mathbf{C}_i$  and modal-specific discriminative information  $\{\mathbf{S}_i^m\}_{m=1}^P$ . Meanwhile,  $\mathbf{C}_i$  is added to the specific modal sequence as a new modal, and then input the new modal sequence  $\{\mathbf{C}_i, \mathbf{S}_i^1, \mathbf{S}_i^2, \dots, \mathbf{S}_i^P\}$  to LSTM structure in order. At  $t$ -th step, we stack all previous hidden output as  $\mathbf{H}_i^t = [h_i^0, h_i^1, \dots, h_i^t]$ , and exploit label correlation  $\mathbf{R}$  based on current modal features. Finally, we make label prediction with stacked output  $\mathbf{H}_i^t$ , label correlation matrix  $\mathbf{R}$  and prediction at  $(t-1)$ -th step.

For multi-modal learning, the goal is to improve performance with heterogeneous modalities or reduce the sample complexity with accumulated multi-modal data [27]. Directed Acyclic Graph (DAG) [20] models an adaptive modalities acquisition system, which learns decision rules that adaptively select modalities for each example as necessary to make a confident prediction. Nevertheless, DAG needs to list all permutations of modal sequence in the training phase. Discriminative Modal Pursuit (DMP) [28] approach is proposed to learn a serialized modal extraction decision methods, which can balance the classification performance and modal feature extraction cost.

What's more, there have been some researches for multi-modal multi-label learning. [2] proposes a new classification framework using the multi-label correlation information to address the problem of simultaneously combining multiple feature modalities and maximum margin classification. [8] performs joint label specific feature selection and take the label correlation matrix as prior knowledge for model training. Considering that label heterogeneity and feature heterogeneity often co-exist, [25] proposes a novel graph-based model for Learning with both Label and Feature heterogeneity ( $L^2F$ ), which imposes the modal consistency by requiring that modal-based classifiers generate similar predictions on the same examples. Multi-Label

Co-Training (MLCT) [23] introduces a predictive reliability measure to select samples, and applies label-wise filtering to confidently communicate labels of selected samples among co-training classifiers. To sufficiently consider the complementary information among multiple modals, LSA-MML [30] is proposed to seek a predictive common representation of multiple modals and the corresponding projection model between the common representation and labels. To further fully extract the complementarity and correlation information effectively, SMISFL [21] jointly learns multiple modal-individual transformations and one sharable transformation. [35] aims to reamin consensus on mutli-modal latent spaces by Hilbert-Schmidt independence criterion during the mapping procedure. However, there is no communication among various modalities. Hence, SIMM [22] is proposed to leverage shared subspace exploitation and view-specific information extraction. Nevertheless, previous approaches rarely consider the label correlation. CS3G approach [29] handles types of interactions between multiple labels, while no interaction between features from different modalities in the model training phase. To make each modality interacts and further reduce modal extraction cost, Multi-modal Classifier Chains (MCC) [34] extends Classifier Chains to exploit label correlations with partial modalities.

### 3 METHODOLOGY

This section mainly gives the detail description of Common and Discriminative Semantic Pursuit (CoDiSP) approach after a preliminary notation explanation.

#### 3.1 Notation

Formally, let  $\mathcal{X} = \mathbb{R}^{d_1} \times \dots \times \mathbb{R}^{d_m} \times \dots \times \mathbb{R}^{d_P}$  be the feature space of  $P$  modalities, where  $d_m (1 \leq m \leq P)$  is the dimensionality of the  $m$ -th modal. Let  $\mathcal{Y} = \{y_k\}_{k=1}^L$  be the label space with  $L$  labels. Given the training dataset with  $N$  data samples  $\mathcal{D} = \{(\mathbf{X}_i, \mathbf{Y}_i)\}_{i=1}^N$ , where  $\mathbf{X}_i = [\mathbf{X}_i^1, \mathbf{X}_i^2, \dots, \mathbf{X}_i^P] \in \mathcal{X}$  is the feature vector and  $\mathbf{Y}_i \subseteq \mathcal{Y}$  is the label vector of the  $i$ -th instance  $\mathbf{X}_i$ . The task of multi-modal multi-label learning is to learn a function  $\mathbf{F} : \mathcal{X} \rightarrow 2^{\mathcal{Y}}$  from  $\mathcal{D}$ , which can assign a set of proper labels for the unseen instance.

#### 3.2 CoDiSP approach

We introduce Common and Discriminative Semantic Pursuit (CoDiSP) approach in detail. As illustrated in Fig. 1, CoDiSP approach contains 4 major components: feature mapping layer, orthogonal layer, LSTM inference layer and label prediction layer.

##### 3.2.1 Feature mapping layer

Considering the dimensionality of different modalities is heterogeneous which is difficult to input to the network, we are supposed to map the original modal feature vectors to the same dimension.

**Common modality mapping** Aiming at exploiting common sub-space representation of all modalities, we concatenate all modalities in  $\mathbf{X}_i$  to formalize a new single modal  $\mathbf{X}_i^0 = [\mathbf{X}_i^1, \mathbf{X}_i^2, \dots, \mathbf{X}_i^P] \in \mathbb{R}^{d_{all}}$ , where  $d_{all} = d_1 + d_2 + \dots + d_P$ . And then we add a dense layer to transform the original common modality to  $d$  dimension common vector according to Eq. 1.

$$\mathbf{C}_i = \text{ReLU}(\mathbf{X}_i^0 \mathbf{U}_0 + \mathbf{b}_0) \quad (1)$$

where  $\mathbf{U}_0 \in \mathbb{R}^{d_{all} \times d}$  is weight vector,  $\mathbf{b}_0 \in \mathbb{R}^{1 \times d}$  is bias vector.

**Specific modality mapping** Each modality in the original feature vector  $\mathbf{X}_i$  can be used to extract specific information. And then we add a dense layer to transform the  $m$ -th original specific modality  $\mathbf{X}_i^m$  with  $d_m$  dimension to  $d$  dimension  $\mathbf{S}_i^m$  by Eq. 2.

$$\mathbf{S}_i^m = \text{ReLU}(\mathbf{X}_i^m \mathbf{U}_m + \mathbf{b}_m), m = 1, \dots, P \quad (2)$$

where  $\mathbf{U}_m \in \mathbb{R}^{d_m \times d}$  is weight vector,  $\mathbf{b}_m \in \mathbb{R}^{1 \times d}$  is bias vector.

As a result, for the  $i$ -th instance  $\mathbf{X}_i$ , feature mapping layer is used to map common and specific modal feature vector into the same  $d$  dimension modal feature vector  $\mathbf{C}_i$  and  $\{\mathbf{S}_i^m\}_{m=1}^P$ , respectively.

##### 3.2.2 Orthogonal layer

Different modalities often includes two main different information: common information shared among all the modalities, and discriminative information of its own specific modality. Therefore, we will introduce our CoDiSP approach based on the following two parts.

**Common semantic pursuit** For traditional single modal multi-label learning approach, concatenating all modalities can obtain better performance than those of the best single modal. As a result, we seek a predictive common representation from concatenated modalities. Furthermore, the corresponding prediction model between common representation and multiple labels can be represented by Eq. 3.

$$\mathbf{f}(\mathbf{C}_i) = \sigma(\mathbf{C}_i \mathbf{W}_C + \mathbf{b}_C) \quad (3)$$

where  $\mathbf{W}_C \in \mathbb{R}^{d \times L}$  is weight vector,  $\mathbf{b}_C \in \mathbb{R}^{1 \times L}$  is the bias vector.

And thus, the common loss function implied in the common features prediction model can be represented as:

$$\mathcal{L}_{comm} = - \sum_{i=1}^{N_b} \sum_{k=1}^L (Y_i^k \log Z_i^k + (1 - Y_i^k) \log(1 - Z_i^k)) \quad (4)$$

where  $Y_i^k$  is the ground-truth of  $\mathbf{X}_i$  on the  $k$ -th label.  $Y_i^k = 1$  if  $k$ -th label is the relevance label, 0 otherwise.  $Z_i^k = \mathbf{f}(\mathbf{C}_i)$  is the prediction with the common modal vector  $\mathbf{C}_i$ , which is extracted from original modal feature vector  $\mathbf{X}_i$ . And  $N_b$  is batch size.

**Discriminative semantic pursuit** Complementary information among different modalities is of great importance. In order to jointly enhance communication between common modality  $\mathbf{X}_i^0$  and other specific modality  $\{\mathbf{X}_i^m\}_{m=1}^P$ , we need to extract specific modal features from  $\{\mathbf{S}_i^m\}_{m=1}^P$ , which eliminate the shared information with common information. We penalize the independence between common modal vector  $\mathbf{C}_i$  and each specific modal vector  $\{\mathbf{S}_i^m\}_{m=1}^P$  with orthogonal loss function:

$$\mathcal{L}_{orth} = \sum_{i=1}^{N_b} \sum_{m=1}^P \|\mathbf{C}_i^T \mathbf{S}_i^m\|_2^2 \quad (5)$$

where  $\|\cdot\|_2$  is the  $L_2$ -norm.  $\mathcal{L}_{orth}$  encourages  $\mathbf{S}_i^m$  extracted from the original  $m$ -th modal vector  $\mathbf{X}_i^m$  to be as discriminative from  $\mathbf{C}_i$  as possible.

##### 3.2.3 LSTM inference layer

After preparing common modal features  $\mathbf{C}_i$  and discriminative modal features  $\{\mathbf{S}_i^m\}_{m=1}^P$ , which are all in the same  $d$  dimension, we input  $\{\mathbf{C}_i, \mathbf{S}_i^1, \dots, \mathbf{S}_i^P\}$  to the network in order. LSTM inference layer has 3 gates as well as 2 states [7]: input gate, forget gate, output gate, cell state and hidden state. At  $t$ -th step, the hidden features of  $\mathbf{X}_i$  in LSTM structure can be represented as  $\mathbf{h}_i^t \in \mathbb{R}^{d_h}$ . To better exploit relationship among different modalities, we stack all the previous hidden outputs as  $\mathbf{H}_i^t = [\mathbf{h}_i^0, \mathbf{h}_i^1, \dots, \mathbf{h}_i^t] \in \mathbb{R}^{(t+1)d_h}$ , where  $d_h$  is the dimension of the hidden layer. All the parameters in LSTM structure are denoted as  $\Psi$ .

##### 3.2.4 Label prediction layer

It is well-known that exploiting label correlations is crucially important in multi-label learning and each modality contains its own specific contribution to the multi-label prediction. In this paper, CoDiSP models label correlations with extracted modal information stored in the memory of LSTM layer.

At the  $t$ -th step, we add a fully connected structure between hidden layer and label prediction layer, which makes label prediction with stacked hidden outputs  $\mathbf{H}_i^t$ . The final label prediction is composed of the prediction of the current modality and the prediction of modality

**Algorithm 1** Training algorithm for CoDiSP approach**Input:**

$\mathcal{D} = \{(\mathbf{X}_i, \mathbf{Y}_i)\}_{i=1}^N$ : Training dataset;  
 $N_b$ : batch size

**Output:**

$\mathbf{F}^P$ : classifier trained with extracted modal sequence

```

1: Concatenate all modalities of  $i$ -th instance  $\mathbf{X}_i$  as  $\mathbf{X}_i^0 = [\mathbf{X}_i^1, \mathbf{X}_i^2, \dots, \mathbf{X}_i^P]$ ,  $i = 1, \dots, N$ 
2: repeat
3:   Randomly select  $N_b$  instances from  $\mathcal{D}$  without replacement
4:   for  $i = 1 : N_b$  do
5:     Map raw common modal  $\mathbf{X}_i^0 \in \mathbb{R}^{d_1 + \dots + d_P}$  to  $\mathbf{C}_i \in \mathbb{R}^d$ 
6:     Map each specific raw modality  $\mathbf{X}_i^m \in \mathbb{R}^{d_m}$  to  $\mathbf{S}_i^m \in \mathbb{R}^d$  with the same dimension,  $m = 1, \dots, P$ 
7:     for  $t = 0 : P$  do
8:       if  $t = 0$  then
9:         Input common modal features  $\mathbf{C}_i$  to LSTM cell
10:      else
11:        Input discriminative modal features  $\mathbf{S}_i^t$  to LSTM cell
12:      end if
13:      Stack hidden output  $\mathbf{H}_i^t = [\mathbf{h}_i^0, \mathbf{h}_i^1, \dots, \mathbf{h}_i^t]$ 
14:      Compute label prediction  $\mathbf{F}^t(\mathbf{H}_i^t)$  with Eq. 6
15:      Compute label loss function  $\mathcal{L}_{pred}^{i,t}$  with Eq. 7
16:    end for
17:  end for
18:  Compute common loss function  $\mathcal{L}_{comm}$  with Eq. 4
19:  Compute orthogonal loss function  $\mathcal{L}_{orth}$  with Eq. 5
20:  Compute overall loss function  $\mathcal{L}$  with Eq. 8
21:  Compute the derivative  $\frac{\partial \mathcal{L}}{\partial \Theta}$ 
22:  Update parameters in  $\Theta$ 
23: until converge
24: return  $\mathbf{F}^P$ 

```

used in the last step. And then we predict multiple labels at  $t$ -th step according to a nonlinear softmax function according to Eq. 6.

$$\mathbf{F}^t(\mathbf{H}_i^t) = \begin{cases} \sigma(\mathbf{H}_i^t \mathbf{W}_L^t + \mathbf{b}_L^t) & t = 0 \\ \sigma(\mathbf{H}_i^t \mathbf{W}_L^t + \mathbf{F}^{t-1}(\mathbf{H}_i^{t-1}) \mathbf{R} + \mathbf{b}_L^t) & t > 0 \end{cases} \quad (6)$$

where  $\mathbf{W}_L^t \in \mathbb{R}^{((t+1)d_h) \times L}$  denotes the fully connected weights between  $\mathbf{H}_i^t$  and label prediction layer,  $\mathbf{b}_L^t \in \mathbb{R}^{1 \times L}$  is the bias vector.

- $\mathbf{H}_i^t \mathbf{W}_L^t$  is similar to BR, which predicts each label independently.
- $\mathbf{F}^{t-1}(\mathbf{H}_i^{t-1}) \mathbf{R}^T$  is the prediction of other labels, in which  $\mathbf{F}^{t-1}(\mathbf{H}_i^{t-1}) \in \mathbb{R}^{1 \times L}$  denotes label prediction at the  $(t-1)$ -th step. Meanwhile, we learn label correlations matrix  $\mathbf{R} \in \mathbb{R}^{L \times L}$ , where  $L$  is the number of labels. The  $k$ -th row and  $j$ -th column of  $\mathbf{R}$  represents the contribution of the  $k$ -th label prediction in  $(t-1)$ -th step to  $j$ -th label, which is denoted as  $R_{kj}$ .

Furthermore, we design binary cross-entropy loss function for final label prediction at  $t$ -th step by Eq. 7.

$$\mathcal{L}_{pred}^{i,t} = - \sum_{k=1}^L (Y_i^k \log \hat{Y}_i^{k,t} + (1 - Y_i^k) \log(1 - \hat{Y}_i^{k,t})) \quad (7)$$

where  $\hat{Y}_i^{k,t}$  the prediction of  $\mathbf{X}_i$  on the  $k$ -th label at  $t$ -th step, predicted by  $\mathbf{F}^t$  in Eq. 6.

Above all, we combine common loss, orthogonal loss and label loss function together to compute the overall loss function  $\mathcal{L}$ :

$$\mathcal{L} = \left( \sum_{i=1}^{N_b} \sum_{t=0}^P \mathcal{L}_{pred}^{i,t} \right) + \alpha \mathcal{L}_{orth} + \beta \mathcal{L}_{comm} \quad (8)$$

where  $\alpha$  and  $\beta$  control the trade-off among different loss terms.

$\Theta = [\mathbf{U}_0, \mathbf{b}_0, \mathbf{U}_m, \mathbf{b}_m, \mathbf{W}_C, \mathbf{b}_C, \mathbf{\Psi}, \mathbf{R}, \mathbf{W}_L^t, \mathbf{b}_L^t]$  denotes all the parameters need to be updated in CoDiSP, where  $m = 1, \dots, P$ ,  $t = 0, \dots, P$ . Then we adopt popular optimization algorithm Adam [13] to update parameters in  $\Theta$  simultaneously. The pseudo code of CoDiSP in the training phase is presented in Algorithm 1.

## 4 EXPERIMENTS

### 4.1 Dataset description

We employ 8 benchmark multi-modal multi-label datasets for performance evaluation, the details are summarized in Table 1.

**Table 1.** Characteristic of the benchmark multi-modal multi-label datasets. #N, #P and #L denote the number of instances, modalities and labels in each dataset, respectively. #D shows the dimensionality of each modality.

Name	#N	#P	#L	#D
<i>Yeast</i>	2417	2	14	[24, 79]
<i>Emotions</i>	593	3	6	[32, 32, 8]
<i>MSRC</i>	591	3	24	[500, 1040, 576]
<i>ML2000</i>	2000	3	5	[500, 1040, 576]
<i>Taobao</i>	2079	4	30	[500, 488, 81, 24]
<i>FCVID</i>	4388	5	28	[400, 400, 400, 400, 400]
<i>Scene</i>	2407	6	6	[49, 49, 49, 49, 49, 49]
<i>MSRA</i>	15000	7	50	[256, 225, 64, 144, 75, 128, 7]

- *Yeast* [32] [22] has two modalities including the genetic expression (70 attributes) and the phylogenetic profile of a gene attributes (24 attributes).
- *Emotions* [19] [34] is a publicly available multi-label dataset with rhythmic attributes and timbre attributes.
- *MSRC* [17] is used for object class recognition. As for each image, there are 3 types of modalities including: BoW, FV and HOG.
- *ML2000* [32] is an image dataset and the modal features are extracted similarly to *MSRC*.
- *Taobao* [29] is used for shopping items classification. Description images of items are crawled from a shopping website, and four types of features, i.e., BoW, Gabor, HOG, HSVHist, are extracted to construct 4 modalities of data. Corresponding categories path of an item provides the label sets.
- *FCVID* [12] is the Fudan-Columbia Video Dataset [11], a subset of 4388 videos with most frequent category names are tested. Each video may come from more than one category and features can be extracted in diverse ways. 5 types of features, namely HOF, HOG, CNN, Trajectory and SIFT are extracted for each video, then PCA is conducted to reduce the dimension of each modal to 400.
- *Scene* [1] [34] is public multi-label dataset with 6 modalities.
- *MSRA* is a subset of a salient object recognition dataset [15], which contains 15000 instances from 50 categories, including 256 RGB color histogram features, 225 dimension block-wise color moments, 64 HSV color histogram, 144 color correlogram, 75 distribution histogram, 128 wavelet features and 7 face features.

**Table 2.** Comparison results (mean  $\pm$  standard deviation) of CoDiSP with compared approaches on benchmark datasets. The best performance for each criterion is bolded.  $\uparrow$  /  $\downarrow$  indicates the larger / smaller the better of the criterion.

Approaches	Hamming Loss $\downarrow$							
	<i>Yeast</i>	<i>Emotions</i>	<i>MSRC</i>	<i>ML2000</i>	<i>Taobao</i>	<i>FCVID</i>	<i>Scene</i>	<i>MSRA</i>
CAMEL(B)	0.193 $\pm$ 0.007	0.218 $\pm$ 0.014	0.059 $\pm$ 0.009	0.089 $\pm$ 0.011	0.032 $\pm$ 0.001	0.018 $\pm$ 0.001	0.144 $\pm$ 0.007	0.046 $\pm$ 0.001
CAMEL(C)	0.189 $\pm$ 0.007	0.207 $\pm$ 0.025	0.106 $\pm$ 0.020	0.098 $\pm$ 0.011	0.033 $\pm$ 0.002	0.020 $\pm$ 0.001	0.076 $\pm$ 0.006	<b>0.045<math>\pm</math>0.001</b>
DMP	0.199 $\pm$ 0.007	0.196 $\pm$ 0.013	0.067 $\pm$ 0.008	0.103 $\pm$ 0.010	0.053 $\pm$ 0.002	0.027 $\pm$ 0.002	0.106 $\pm$ 0.007	0.046 $\pm$ 0.001
CS3G	0.255 $\pm$ 0.008	0.290 $\pm$ 0.013	0.077 $\pm$ 0.008	0.119 $\pm$ 0.010	0.064 $\pm$ 0.003	0.020 $\pm$ 0.001	0.217 $\pm$ 0.013	0.050 $\pm$ 0.001
MCC	0.213 $\pm$ 0.009	0.214 $\pm$ 0.023	0.068 $\pm$ 0.008	0.105 $\pm$ 0.012	0.054 $\pm$ 0.002	0.027 $\pm$ 0.001	0.102 $\pm$ 0.010	0.048 $\pm$ 0.001
CoDiSP	<b>0.189<math>\pm</math>0.009</b>	<b>0.173<math>\pm</math>0.025</b>	<b>0.050<math>\pm</math>0.008</b>	<b>0.086<math>\pm</math>0.011</b>	<b>0.027<math>\pm</math>0.002</b>	<b>0.012<math>\pm</math>0.001</b>	<b>0.064<math>\pm</math>0.009</b>	0.046 $\pm$ 0.001
Approaches	Ranking Loss $\downarrow$							
	<i>Yeast</i>	<i>Emotions</i>	<i>MSRC</i>	<i>ML2000</i>	<i>Taobao</i>	<i>FCVID</i>	<i>Scene</i>	<i>MSRA</i>
CAMEL(B)	0.162 $\pm$ 0.011	0.176 $\pm$ 0.023	<b>0.033<math>\pm</math>0.009</b>	<b>0.063<math>\pm</math>0.011</b>	0.151 $\pm$ 0.014	0.027 $\pm$ 0.006	0.158 $\pm$ 0.018	0.159 $\pm$ 0.010
CAMEL(C)	0.163 $\pm$ 0.012	0.179 $\pm$ 0.038	0.153 $\pm$ 0.074	0.067 $\pm$ 0.010	0.159 $\pm$ 0.012	0.031 $\pm$ 0.005	0.057 $\pm$ 0.011	0.154 $\pm$ 0.009
DMP	0.200 $\pm$ 0.010	0.150 $\pm$ 0.021	0.047 $\pm$ 0.010	0.074 $\pm$ 0.010	0.238 $\pm$ 0.020	0.052 $\pm$ 0.007	0.088 $\pm$ 0.007	0.204 $\pm$ 0.005
CS3G	0.211 $\pm$ 0.012	0.225 $\pm$ 0.018	0.043 $\pm$ 0.010	0.092 $\pm$ 0.013	0.171 $\pm$ 0.009	<b>0.027<math>\pm</math>0.003</b>	0.289 $\pm$ 0.030	0.140 $\pm$ 0.007
MCC	0.224 $\pm$ 0.016	0.178 $\pm$ 0.029	0.054 $\pm$ 0.011	0.082 $\pm$ 0.011	0.235 $\pm$ 0.016	0.052 $\pm$ 0.005	0.101 $\pm$ 0.011	0.195 $\pm$ 0.005
CoDiSP	<b>0.159<math>\pm</math>0.010</b>	<b>0.136<math>\pm</math>0.022</b>	0.034 $\pm$ 0.012	0.064 $\pm$ 0.013	<b>0.098<math>\pm</math>0.013</b>	0.029 $\pm$ 0.004	<b>0.052<math>\pm</math>0.007</b>	<b>0.120<math>\pm</math>0.005</b>
Approaches	Subset Accuracy $\uparrow$							
	<i>Yeast</i>	<i>Emotions</i>	<i>MSRC</i>	<i>ML2000</i>	<i>Taobao</i>	<i>FCVID</i>	<i>Scene</i>	<i>MSRA</i>
CAMEL(B)	0.173 $\pm$ 0.022	0.248 $\pm$ 0.034	0.322 $\pm$ 0.058	0.655 $\pm$ 0.040	0.150 $\pm$ 0.020	0.534 $\pm$ 0.030	0.300 $\pm$ 0.033	0.057 $\pm$ 0.010
CAMEL(C)	0.201 $\pm$ 0.018	0.272 $\pm$ 0.048	0.070 $\pm$ 0.073	0.633 $\pm$ 0.036	0.238 $\pm$ 0.046	0.485 $\pm$ 0.019	0.646 $\pm$ 0.024	0.066 $\pm$ 0.010
DMP	0.151 $\pm$ 0.022	0.284 $\pm$ 0.058	0.291 $\pm$ 0.033	0.617 $\pm$ 0.031	0.218 $\pm$ 0.017	0.520 $\pm$ 0.033	0.531 $\pm$ 0.032	0.053 $\pm$ 0.005
CS3G	0.048 $\pm$ 0.009	0.165 $\pm$ 0.041	0.205 $\pm$ 0.043	0.564 $\pm$ 0.033	0.104 $\pm$ 0.017	0.571 $\pm$ 0.018	0.327 $\pm$ 0.037	0.061 $\pm$ 0.009
MCC	0.190 $\pm$ 0.034	0.299 $\pm$ 0.037	0.345 $\pm$ 0.088	0.662 $\pm$ 0.032	0.213 $\pm$ 0.023	0.522 $\pm$ 0.024	0.662 $\pm$ 0.038	0.076 $\pm$ 0.003
CoDiSP	<b>0.222<math>\pm</math>0.030</b>	<b>0.403<math>\pm</math>0.069</b>	<b>0.470<math>\pm</math>0.058</b>	<b>0.745<math>\pm</math>0.034</b>	<b>0.527<math>\pm</math>0.040</b>	<b>0.783<math>\pm</math>0.026</b>	<b>0.777<math>\pm</math>0.030</b>	<b>0.148<math>\pm</math>0.010</b>
Approaches	Macro F1 $\uparrow$							
	<i>Yeast</i>	<i>Emotions</i>	<i>MSRC</i>	<i>ML2000</i>	<i>Taobao</i>	<i>FCVID</i>	<i>Scene</i>	<i>MSRA</i>
CAMEL(B)	0.392 $\pm$ 0.025	0.590 $\pm$ 0.039	0.680 $\pm$ 0.043	0.804 $\pm$ 0.025	0.034 $\pm$ 0.008	0.697 $\pm$ 0.022	0.458 $\pm$ 0.033	0.069 $\pm$ 0.004
CAMEL(C)	0.458 $\pm$ 0.020	0.615 $\pm$ 0.058	0.208 $\pm$ 0.041	0.781 $\pm$ 0.023	0.182 $\pm$ 0.028	0.659 $\pm$ 0.017	0.772 $\pm$ 0.021	0.079 $\pm$ 0.002
DMP	0.322 $\pm$ 0.012	0.628 $\pm$ 0.032	0.643 $\pm$ 0.045	0.781 $\pm$ 0.020	0.231 $\pm$ 0.015	0.660 $\pm$ 0.023	0.689 $\pm$ 0.020	0.054 $\pm$ 0.002
CS3G	0.217 $\pm$ 0.017	0.461 $\pm$ 0.023	0.506 $\pm$ 0.040	0.738 $\pm$ 0.021	0.125 $\pm$ 0.015	0.693 $\pm$ 0.014	0.282 $\pm$ 0.034	0.041 $\pm$ 0.007
MCC	0.340 $\pm$ 0.018	0.621 $\pm$ 0.039	0.652 $\pm$ 0.024	0.780 $\pm$ 0.027	0.222 $\pm$ 0.026	0.667 $\pm$ 0.009	0.718 $\pm$ 0.026	0.073 $\pm$ 0.004
CoDiSP	<b>0.472<math>\pm</math>0.022</b>	<b>0.714<math>\pm</math>0.043</b>	<b>0.755<math>\pm</math>0.047</b>	<b>0.828<math>\pm</math>0.021</b>	<b>0.398<math>\pm</math>0.028</b>	<b>0.834<math>\pm</math>0.020</b>	<b>0.830<math>\pm</math>0.022</b>	<b>0.194<math>\pm</math>0.013</b>
Approaches	Example F1 $\uparrow$							
	<i>Yeast</i>	<i>Emotions</i>	<i>MSRC</i>	<i>ML2000</i>	<i>Taobao</i>	<i>FCVID</i>	<i>Scene</i>	<i>MSRA</i>
CAMEL(B)	0.618 $\pm$ 0.014	0.534 $\pm$ 0.039	0.805 $\pm$ 0.032	0.769 $\pm$ 0.031	0.054 $\pm$ 0.012	0.551 $\pm$ 0.031	0.341 $\pm$ 0.035	0.232 $\pm$ 0.010
CAMEL(C)	0.628 $\pm$ 0.013	0.581 $\pm$ 0.049	0.618 $\pm$ 0.074	0.739 $\pm$ 0.029	0.266 $\pm$ 0.044	0.503 $\pm$ 0.017	0.695 $\pm$ 0.027	0.248 $\pm$ 0.006
DMP	0.609 $\pm$ 0.017	0.605 $\pm$ 0.037	0.789 $\pm$ 0.028	0.753 $\pm$ 0.026	0.336 $\pm$ 0.014	0.642 $\pm$ 0.026	0.620 $\pm$ 0.026	0.216 $\pm$ 0.007
CS3G	0.549 $\pm$ 0.011	0.538 $\pm$ 0.028	0.729 $\pm$ 0.024	0.697 $\pm$ 0.029	0.323 $\pm$ 0.025	0.673 $\pm$ 0.022	0.366 $\pm$ 0.036	0.273 $\pm$ 0.011
MCC	0.585 $\pm$ 0.025	0.624 $\pm$ 0.043	0.798 $\pm$ 0.024	0.787 $\pm$ 0.022	0.333 $\pm$ 0.026	0.647 $\pm$ 0.019	0.713 $\pm$ 0.029	0.266 $\pm$ 0.005
CoDiSP	<b>0.629<math>\pm</math>0.018</b>	<b>0.690<math>\pm</math>0.050</b>	<b>0.850<math>\pm</math>0.026</b>	<b>0.827<math>\pm</math>0.022</b>	<b>0.541<math>\pm</math>0.039</b>	<b>0.796<math>\pm</math>0.026</b>	<b>0.818<math>\pm</math>0.023</b>	<b>0.341<math>\pm</math>0.013</b>
Approaches	Micro F1 $\uparrow$							
	<i>Yeast</i>	<i>Emotions</i>	<i>MSRC</i>	<i>ML2000</i>	<i>Taobao</i>	<i>FCVID</i>	<i>Scene</i>	<i>MSRA</i>
CAMEL(B)	0.642 $\pm$ 0.012	0.607 $\pm$ 0.038	0.814 $\pm$ 0.029	0.806 $\pm$ 0.024	0.101 $\pm$ 0.023	0.695 $\pm$ 0.024	0.457 $\pm$ 0.037	0.329 $\pm$ 0.014
CAMEL(C)	0.658 $\pm$ 0.012	0.637 $\pm$ 0.048	0.629 $\pm$ 0.071	0.783 $\pm$ 0.024	0.373 $\pm$ 0.053	0.658 $\pm$ 0.016	0.763 $\pm$ 0.019	0.349 $\pm$ 0.010
DMP	0.632 $\pm$ 0.014	0.658 $\pm$ 0.027	0.796 $\pm$ 0.025	0.783 $\pm$ 0.021	0.359 $\pm$ 0.015	0.658 $\pm$ 0.023	0.683 $\pm$ 0.019	0.311 $\pm$ 0.009
CS3G	0.566 $\pm$ 0.011	0.574 $\pm$ 0.031	0.744 $\pm$ 0.026	0.744 $\pm$ 0.021	0.332 $\pm$ 0.026	0.720 $\pm$ 0.015	0.372 $\pm$ 0.036	0.324 $\pm$ 0.018
MCC	0.615 $\pm$ 0.021	0.653 $\pm$ 0.041	0.799 $\pm$ 0.023	0.784 $\pm$ 0.024	0.354 $\pm$ 0.021	0.663 $\pm$ 0.014	0.709 $\pm$ 0.027	0.359 $\pm$ 0.006
CoDiSP	<b>0.662<math>\pm</math>0.015</b>	<b>0.719<math>\pm</math>0.044</b>	<b>0.851<math>\pm</math>0.025</b>	<b>0.826<math>\pm</math>0.021</b>	<b>0.591<math>\pm</math>0.041</b>	<b>0.828<math>\pm</math>0.021</b>	<b>0.820<math>\pm</math>0.023</b>	<b>0.430<math>\pm</math>0.012</b>

**Table 3.** Comparison results (mean  $\pm$  standard deviation) of CoDiSP-NC and CoDiSP, where CoDiSP-NC denotes the model using only discriminative information  $\{S^m\}_{m=1}^P$  of each modality in the CoDiSP network. The best performance for each criterion is bolded.  $\uparrow$  /  $\downarrow$  indicates the larger / smaller the better of the criterion.

Datasets	Approaches	Evaluation Metrics					
		Hamming Loss $\downarrow$	Ranking Loss $\downarrow$	Subset Accuracy $\uparrow$	Macro F1 $\uparrow$	Example F1 $\uparrow$	Micro F1 $\uparrow$
Yeast	CoDiSP-NC	0.195 $\pm$ 0.008	0.164 $\pm$ 0.009	0.214 $\pm$ 0.025	0.461 $\pm$ 0.016	0.617 $\pm$ 0.015	0.649 $\pm$ 0.013
	CoDiSP	<b>0.189<math>\pm</math>0.009</b>	<b>0.159<math>\pm</math>0.010</b>	<b>0.222<math>\pm</math>0.030</b>	<b>0.472<math>\pm</math>0.022</b>	<b>0.629<math>\pm</math>0.018</b>	<b>0.662<math>\pm</math>0.015</b>
Emotions	CoDiSP-NC	0.174 $\pm$ 0.021	0.136 $\pm$ 0.019	0.393 $\pm$ 0.065	0.707 $\pm$ 0.038	0.664 $\pm$ 0.039	0.710 $\pm$ 0.040
	CoDiSP	<b>0.173<math>\pm</math>0.025</b>	<b>0.136<math>\pm</math>0.022</b>	<b>0.403<math>\pm</math>0.069</b>	<b>0.714<math>\pm</math>0.043</b>	<b>0.690<math>\pm</math>0.050</b>	<b>0.719<math>\pm</math>0.044</b>
MSRC	CoDiSP-NC	0.055 $\pm$ 0.012	<b>0.033<math>\pm</math>0.012</b>	0.455 $\pm$ 0.063	0.735 $\pm$ 0.066	0.838 $\pm$ 0.039	0.840 $\pm$ 0.037
	CoDiSP	<b>0.050<math>\pm</math>0.008</b>	0.034 $\pm$ 0.012	<b>0.470<math>\pm</math>0.058</b>	<b>0.755<math>\pm</math>0.047</b>	<b>0.850<math>\pm</math>0.026</b>	<b>0.851<math>\pm</math>0.025</b>
ML2000	CoDiSP-NC	0.091 $\pm$ 0.012	0.072 $\pm$ 0.013	0.731 $\pm$ 0.032	0.816 $\pm$ 0.022	0.815 $\pm$ 0.024	0.813 $\pm$ 0.022
	CoDiSP	<b>0.086<math>\pm</math>0.011</b>	<b>0.064<math>\pm</math>0.013</b>	<b>0.745<math>\pm</math>0.034</b>	<b>0.828<math>\pm</math>0.021</b>	<b>0.827<math>\pm</math>0.022</b>	<b>0.826<math>\pm</math>0.021</b>
Taobao	CoDiSP-NC	0.029 $\pm$ 0.002	0.113 $\pm$ 0.015	0.471 $\pm$ 0.034	0.346 $\pm$ 0.028	0.473 $\pm$ 0.035	0.540 $\pm$ 0.041
	CoDiSP	<b>0.027<math>\pm</math>0.002</b>	<b>0.098<math>\pm</math>0.013</b>	<b>0.527<math>\pm</math>0.040</b>	<b>0.398<math>\pm</math>0.028</b>	<b>0.541<math>\pm</math>0.039</b>	<b>0.591<math>\pm</math>0.041</b>
FCVID	CoDiSP-NC	0.015 $\pm$ 0.001	0.052 $\pm$ 0.007	0.708 $\pm$ 0.024	0.780 $\pm$ 0.014	0.718 $\pm$ 0.021	0.777 $\pm$ 0.015
	CoDiSP	<b>0.012<math>\pm</math>0.001</b>	<b>0.029<math>\pm</math>0.004</b>	<b>0.783<math>\pm</math>0.026</b>	<b>0.834<math>\pm</math>0.020</b>	<b>0.796<math>\pm</math>0.026</b>	<b>0.828<math>\pm</math>0.021</b>
Scene	CoDiSP-NC	0.069 $\pm$ 0.008	0.055 $\pm$ 0.010	0.757 $\pm$ 0.025	0.814 $\pm$ 0.019	0.792 $\pm$ 0.022	0.803 $\pm$ 0.021
	CoDiSP	<b>0.064<math>\pm</math>0.009</b>	<b>0.052<math>\pm</math>0.007</b>	<b>0.777<math>\pm</math>0.030</b>	<b>0.830<math>\pm</math>0.022</b>	<b>0.818<math>\pm</math>0.023</b>	<b>0.820<math>\pm</math>0.023</b>
MSRA	CoDiSP-NC	0.046 $\pm$ 0.001	0.132 $\pm$ 0.004	0.115 $\pm$ 0.007	0.101 $\pm$ 0.007	0.259 $\pm$ 0.016	0.353 $\pm$ 0.015
	CoDiSP	<b>0.046<math>\pm</math>0.001</b>	<b>0.120<math>\pm</math>0.005</b>	<b>0.148<math>\pm</math>0.010</b>	<b>0.194<math>\pm</math>0.013</b>	<b>0.341<math>\pm</math>0.013</b>	<b>0.430<math>\pm</math>0.012</b>

## 4.2 Compared approaches

Considering CoDiSP is related to multi-modal multi-label learning, the performance of CoDiSP is compared against 5 approaches, including a state-of-the-art multi-label learning approach with two types of feature inputs and 3 multi-modal multi-label approaches.

- CAMEL(B) & CAMEL(C): CAMEL [3] is a novel multi-label learning approach that aims to explicitly account for the correlated predictions of labels while training the desired model simultaneously. CAMEL(B) stands for the best performance obtained from the best single modality. CAMEL(C) stands for concatenating all the modalities as a single modal input.
- DMP [28]: A multi-modal learning approach, which predicts the label information and decide the modalities to be extracted simultaneously. Here we treat each label independently with DMP.
- CS3G [29]: A multi-modal multi-label approach utilizes multi-modal information in a privacy-preserving style to deal with multi-label tasks, which treats each modality unequally and has the ability to extract the most useful modal features for final prediction.
- MCC [34]: A multi-modal multi-label approach that makes great use of modalities, and can make a convince prediction with many instead of all modalities.

## 4.3 Evaluation metrics

To evaluate the performance of CoDiSP compared with other multi-modal multi-label learning approaches, we adopt 6 widely-used multi-label evaluation metrics, including Hamming Loss, Ranking Loss, Subset Accuracy, Macro F1, Example F1 and Micro F1 [33]. For Hamming Loss and Ranking Loss, smaller value indicates better performance, while larger value of the other 4 evaluation metrics means better performance. And all the employed evaluation metrics vary within the interval  $[0, 1]$ .

## 4.4 Experimental results

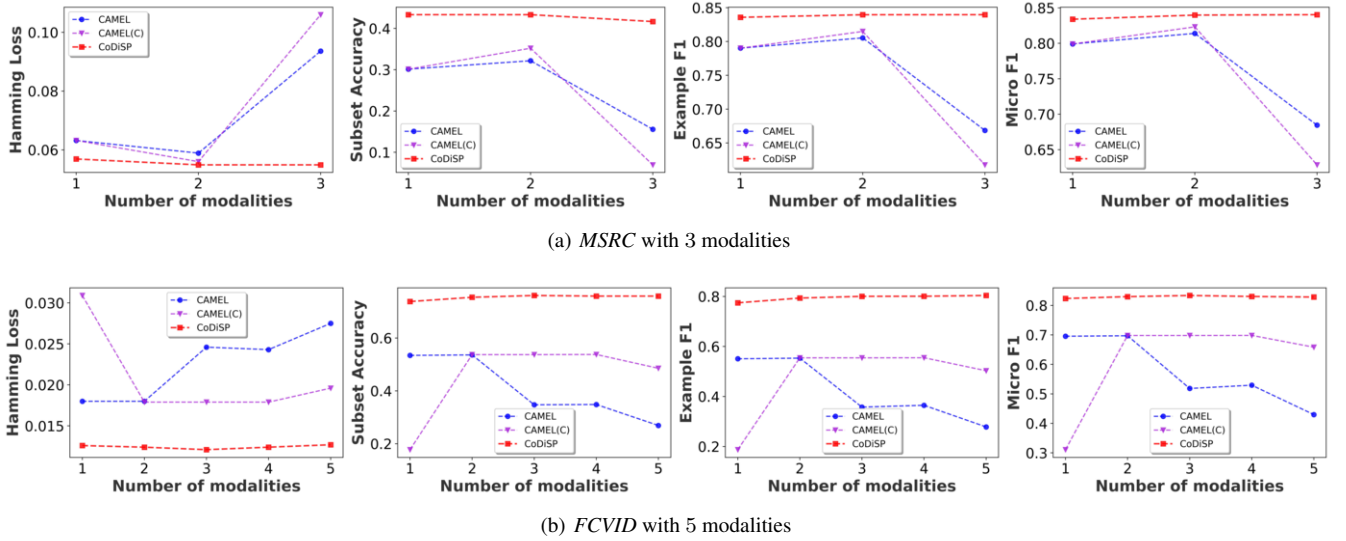
### 4.4.1 Comparison with state-of-the-arts

For all these approaches, we report the best results of the optimal parameters in terms of classification performance. Meanwhile, we perform 10-fold cross validation (CV) and take the average value of the results in the end. We set trade-off parameter  $\alpha = 0.1$ ,  $\beta = 100$  and batch size  $N_b = 64$ . Furthermore, we drop out 40% modal features at each step to avoid over-fitting [18].

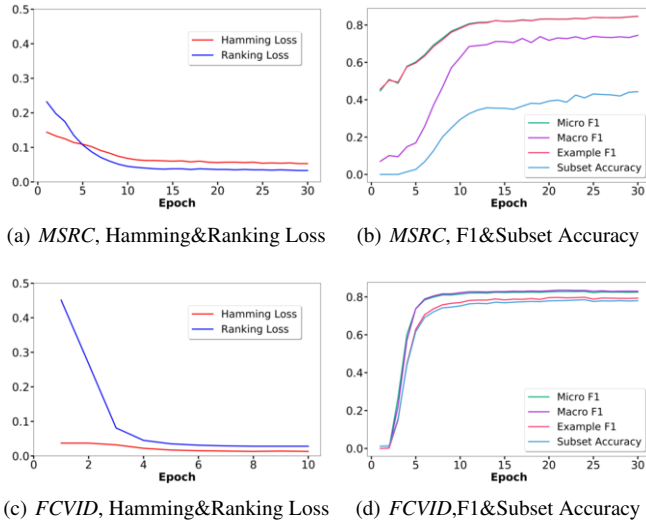
Table 2 demonstrates the comparison results of different approaches on benchmark datasets. Based on the experimental results, several observations are obtained as follows: 1) CoDiSP approach achieves the best performance on all benchmark datasets compared with other state-of-the-art approaches, which demonstrates the highly competitive performance of our proposed approach in multi-modal multi-label learning problem. 2) Concatenating all modalities as a single modality may not always achieve better performance than those of best single modal, which indicates the necessity of taking both concatenated modalities and single modality into consideration.

### 4.4.2 Common Semantic Analysis

In order to validate the effectiveness of common information extraction, we keep the basic structure of our proposed CoDiSP model and remove the common features  $C$  from the modal sequence, which is denoted as CoDiSP-NC. As shown in Table 3, CoDiSP performs better than CoDiSP-NC, which shows the significance of common information exploitation of all the modalities, i.e., it is not good enough to merely fuse semantic of specific modalities. Even though, CoDiSP-NC performs better than other state-of-the-art approaches shown in Table 2, which shows the effectiveness of exploiting discriminative information from each specific modality.



**Figure 2.** Performance of CAMEL, CAMEL(C) and CoDiSP with increase of the modality on the *MSRC* and *FCVID* dataset. With the emergence of new modality, CAMEL adopts the new modality as input, while CAMEL(C) concatenates new modality with previous modalities as input.



**Figure 3.** Convergence analysis of CoDiSP approach on the *MSRC* and *FCVID* dataset.

#### 4.4.3 Modal dependency analysis

To further evaluate the effectiveness of CoDiSP in sequentially learning modal dependency without pre-determining the modal order, we provide comparison for the state-of-the-art multi-label approach CAMEL in terms of Hamming Loss, Subset Accuracy, Example F1 and Micro F1. Based on the comparison results in Fig. 2, we observe that: 1) CAMEL performs differently with different modalities, so it is significant to exploit discriminative information of each specific modality. 2) With the emergence of new modality, the performance of CAMEL(C) is getting worse on *MSRC* dataset, while the performance of CAMEL(C) is getting better on *FCVID* dataset. In contrast, the curve of CoDiSP tends to be more stable. This phenomenon shows that concatenating modalities is not a wise decision,

and extraction of the common features of all modalities with orthogonal constraint is effective.

#### 4.4.4 Convergence analysis

We conduct convergence experiments to validate the convergence of CoDiSP, due to the page limit, we only give the convergence results on two datasets: *MSRC* and *FCVID*. As shown in Fig. 3, CoDiSP approach can converge fast within a small number of epochs.

## 5 CONCLUSION

Rapid development of data collection techniques has spawned the multi-modal multi-label learning, while the modalities extracted from different channels are inconsistent. In this paper, a novel multi-modal multi-label learning approach named Common and Discriminative Semantic Pursuit (CoDiSP) is proposed to exploit consistent and complementary information of multiple modalities, which encodes the commonality and discrimination of different modalities in latent semantic space. And then, common information is added to the specific modal sequence as a new modal. What's more, CoDiSP explores modal relationship and hidden label correlations when inputting into the LSTM network with features of each modality in the new modal sequence. Experiments on 8 benchmark datasets validate the effectiveness of our proposed CoDiSP approach compared with other state-of-the-art approaches. Furthermore, we conduct extensive experiments to analyze common semantic pursuit, modal dependency and convergence. In the future, how to transfer the learned model between different domains is an interesting work.

## ACKNOWLEDGEMENTS

This paper is supported by the National Key Research and Development Program of China (Grant No. 2018YFB1403400), the National Natural Science Foundation of China (Grant No. 61876080), the Key Research and Development Program of Jiangsu (Grant No. BE2019105), the Collaborative Innovation Center of Novel Software Technology and Industrialization at Nanjing University.

## REFERENCES

- [1] Matthew R Boutell, Jiebo Luo, Xipeng Shen, and Christopher M Brown, 'Learning multi-label scene classification', *Pattern recognition*, **37**(9), 1757–1771, (2004).
- [2] Zheng Fang and Zhongfei Zhang, 'Simultaneously combining multi-view multi-label learning with maximum margin classification', in *2012 IEEE 12th International Conference on Data Mining*, pp. 864–869. IEEE, (2012).
- [3] Lei Feng, Bo An, and Shuo He, 'Collaboration based multi-label learning', in *Thirty-Third AAAI Conference on Artificial Intelligence*, pp. 3550–3557, (2019).
- [4] Johannes Fürnkranz, Eyke Hüllermeier, Eneldo Loza Mencía, and Klaus Brinker, 'Multilabel classification via calibrated label ranking', *Machine learning*, **73**(2), 133–153, (2008).
- [5] Eva Gibaja and Sebastián Ventura, 'A tutorial on multilabel learning', *ACM Computing Surveys (CSUR)*, **47**(3), 52, (2015).
- [6] Siddharth Gopal and Yiming Yang, 'Multilabel classification with meta-level features', in *Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval*, pp. 315–322. ACM, (2010).
- [7] Sepp Hochreiter and Jürgen Schmidhuber, 'Long short-term memory', *Neural computation*, **9**(8), 1735–1780, (1997).
- [8] Jun Huang, Guorong Li, Qingming Huang, and Xindong Wu, 'Learning label specific features for multi-label classification', in *2015 IEEE International Conference on Data Mining*, pp. 181–190. IEEE, (2015).
- [9] Jun Huang, Guorong Li, Qingming Huang, and Xindong Wu, 'Joint feature selection and classification for multilabel learning', *IEEE transactions on cybernetics*, **48**(3), 876–889, (2017).
- [10] Sheng-Jun Huang, Yang Yu, and Zhi-Hua Zhou, 'Multi-label hypothesis reuse', in *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 525–533. ACM, (2012).
- [11] Sheng-Jun Huang and Zhi-Hua Zhou, 'Multi-label learning by exploiting label correlations locally', in *Twenty-sixth AAAI conference on artificial intelligence*, (2012).
- [12] Yu-Gang Jiang, Zuxuan Wu, Jun Wang, Xiangyang Xue, and Shih-Fu Chang, 'Exploiting feature and class relationships in video categorization with regularized deep neural networks', *IEEE Transactions on Pattern Analysis & Machine Intelligence*, **40**(2), 352–364, (2018).
- [13] Diederik P Kingma and Jimmy Ba, 'Adam: A method for stochastic optimization', in *Proceedings of the 3rd International Conference on Learning Representations*, (2015).
- [14] Ximing Li, Jihong Ouyang, and Xiaotang Zhou, 'Supervised topic models for multi-label classification', *Neurocomputing*, **149**, 811–819, (2015).
- [15] Tie Liu, Zejian Yuan, Jian Sun, Jingdong Wang, Nanning Zheng, Xiaou Tang, and Heung-Yeung Shum, 'Learning to detect a salient object', *IEEE Transactions on Pattern analysis and machine intelligence*, **33**(2), 353–367, (2010).
- [16] Jesse Read, Bernhard Pfahringer, Geoff Holmes, and Eibe Frank, 'Classifier chains for multi-label classification', *Machine learning*, **85**(3), 333, (2011).
- [17] Florian Schroff, Antonio Criminisi, and Andrew Zisserman, 'Harvesting image databases from the web', *IEEE transactions on pattern analysis and machine intelligence*, **33**(4), 754–766, (2010).
- [18] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov, 'Dropout: a simple way to prevent neural networks from overfitting', *The journal of machine learning research*, **15**(1), 1929–1958, (2014).
- [19] Konstantinos Trohidis, Grigorios Tsoumakas, George Kalliris, and Ioannis P Vlahavas, 'Multi-label classification of music into emotions.', in *ISMIR*, volume 8, pp. 325–330, (2008).
- [20] Joseph Wang, Kirill Trapeznikov, and Venkatesh Saligrama, 'Efficient learning by directed acyclic graph for resource constrained prediction', in *Advances in Neural Information Processing Systems*, pp. 2152–2160, (2015).
- [21] Fei Wu, Xiao-Yuan Jing, Jun Zhou, Yimu Ji, Chao Lan, Qinghua Huang, and Ruchuan Wang, 'Semi-supervised multi-view individual and sharable feature learning for webpage classification', in *The World Wide Web Conference*, pp. 3349–3355. ACM, (2019).
- [22] Xuan Wu, Qing-Guo Chen, Yao Hu, Dengbao Wang, Xiaodong Chang, Xiaobo Wang, and Min-Ling Zhang, 'Multi-view multi-label learning with view-specific information extraction', in *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pp. 3884–3890. AAAI Press, (2019).
- [23] Yuying Xing, Guoxian Yu, Carlotta Domeniconi, Jun Wang, and Zili Zhang, 'Multi-label co-training', in *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pp. 2882–2888. AAAI Press, (2018).
- [24] Hao Yang, Joey Tianyi Zhou, Yu Zhang, Bin-Bin Gao, Jianxin Wu, and Jianfei Cai, 'Exploit bounding box annotations for multi-label object recognition', in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 280–288, (2016).
- [25] Pei Yang, Hongxia Yang, Haoda Fu, Dawei Zhou, Jieping Ye, Theodoros Lappas, and Jingrui He, 'Jointly modeling label and feature heterogeneity in medical informatics', *ACM Transactions on Knowledge Discovery from Data (TKDD)*, **10**(4), 39, (2016).
- [26] Yang Yang, Ke-Tao Wang, De-Chuan Zhan, Hui Xiong, and Yuan Jiang, 'Comprehensive semi-supervised multi-modal learning', in *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pp. 4092–4098. AAAI Press, (2019).
- [27] Yang Yang, Yi-Feng Wu, De-Chuan Zhan, Zhi-Bin Liu, and Yuan Jiang, 'Complex object classification: A multi-modal multi-instance multi-label deep network with optimal transport', in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 2594–2603. ACM, (2018).
- [28] Yang Yang, De-Chuan Zhan, Ying Fan, and Yuan Jiang, 'Instance specific discriminative modal pursuit: A serialized approach', in *Asian Conference on Machine Learning*, pp. 65–80, (2017).
- [29] Han-Jia Ye, De-Chuan Zhan, Xiaolin Li, Zhen-Chuan Huang, and Yuan Jiang, 'College student scholarships and subsidies granting: A multi-modal multi-label approach', in *2016 IEEE 16th International Conference on Data Mining (ICDM)*, pp. 559–568. IEEE, (2016).
- [30] Changqing Zhang, Ziwei Yu, Qinghua Hu, Pengfei Zhu, Xinwang Liu, and Xiaobo Wang, 'Latent semantic aware multi-view multi-label classification', in *Thirty-Second AAAI Conference on Artificial Intelligence*, (2018).
- [31] Min-Ling Zhang and Zhi-Hua Zhou, 'Multilabel neural networks with applications to functional genomics and text categorization', *IEEE transactions on Knowledge and Data Engineering*, **18**(10), 1338–1351, (2006).
- [32] Min Ling Zhang and Zhi Hua Zhou, 'Ml-knn: A lazy learning approach to multi-label learning', *Pattern Recognition*, **40**(7), 2038–2048, (2007).
- [33] Min-Ling Zhang and Zhi-Hua Zhou, 'A review on multi-label learning algorithms', *IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING*, **26**(8), 1819, (2014).
- [34] Yi Zhang, Cheng Zeng, Hao Cheng, Chongjun Wang, and Lei Zhang, 'Many could be better than all: A novel instance-oriented algorithm for multi-modal multi-label problem', in *2019 IEEE International Conference on Multimedia and Expo (ICME)*, pp. 838–843. IEEE, (2019).
- [35] Pengfei Zhu, Qi Hu, Qinghua Hu, Changqing Zhang, and Zhizhao Feng, 'Multi-view label embedding', *Pattern Recognition*, **84**, 126–135, (2018).