

# Rethinking Modal-oriented Label Correlations for Multi-modal Multi-label Learning

Yi Zhang, Jundong Shen, Zhecheng Zhang, Lei Zhang, Chongjun Wang\*

National Key Laboratory for Novel Software Technology

Nanjing University, Nanjing 210023, China

{njuzhangyi, jdshen, zzc}@smail.nju.edu.cn, {zhangl, chjwang}@nju.edu.cn

**Abstract**—Multi-modal multi-label learning provides a fundamental framework for complex objects, which can be represented with multiple modalities and annotated with multiple labels simultaneously. Different modalities can usually provide complementary information, which may lead to improved performance. What's more, exploiting label correlations is crucially important to multi-label learning. However, most existing multi-label learning approaches do not sufficiently consider the complementary information among different modalities. In this paper, we propose a novel end-to-end deep learning framework named Rethinking Modal-oriented Label Correlations (RMLC), which sequentially polish the label prediction with each individual modality. In order to explicitly account for the correlated prediction of multiple labels, RMLC leverages an efficient sequential modal-based exploration to rethink label correlations. The final prediction of each label involves the collaboration between modal-specific prediction and the prediction of other labels based on cross-modal interaction. Comprehensive experiments on benchmark datasets validate the effectiveness and competitiveness of the proposed RMLC approach.

**Index Terms**—multi-modal, multi-label, label correlations, modal-specific, cross-modal

## I. INTRODUCTION

With the fast development of data collection techniques, objects are often characterized by features from different data channels, i.e., multi-modal feature representations. For example, a news webpage can be represented with two heterogeneous modalities: text and image; an image can be described using different features, such as texture descriptors, shape descriptors, color descriptors, surrounding texts, and so on [1].

Although multi-modal (or multi-view) learning approaches have been developed and paid more attention, most previous studies assume that each object is annotated with a single label. Nevertheless, in real-world applications, each object may have multiple semantic meanings. For instance, a webpage may be tagged with multiple labels: economics, sports and culture.

As a result, Multi-Modal Multi-Label (MMML) learning serves an important framework to solve complex objects with multiple modalities and multiple labels. For example, in video annotation, a film can be represented from multiple channels including text, audio, picture and frame, meanwhile it can be annotated with *superhero movie* (type), *Marvel Studios* (producer), *America* (country) and *Anthony Russo and Joseph Russo* (directors). The major challenge of MMML learning lies in how to jointly model the multiple types of heterogeneities

in a mutually beneficial way. For one thing, the representations of various modalities are quite different from each other and it is a challenge to fuse the multiple modalities directly with large discrepancy. Previous approaches do not explicitly account for the distinctive information hidden in each specific modality. For another, how to exploit label correlations in an effectiveness way is also a challenging issue. Since each specific modality captures a specific property of data, it is impossible for one modality to comprehensively characterize all the relevant labels.

In this paper, aiming at simultaneously exploit the modal correlations and label correlations in a mutually benefit way, we proposed a novel MMML learning approach named Rethinking Modal-oriented Label Correlations (RMLC). On the one hand, RMLC models the cross-modal interaction (i.e., rethinking process) in a LSTM network [2], which captures the complementary patterns among multiple modalities. On the other hand, we input each specific modality to the RMLC network step by step, in which we make label prediction and exploit label correlations simultaneously. For each label, the final prediction consists of two aspects: (1) its own prediction based on each modal-specific information (2) the prediction of other labels based on cross-modal interaction.

The main contributions are summarized as follows.

- We propose a novel end-to-end deep network structure named Rethinking Modal-oriented Label Correlations (RMLC) for multi-modal multi-label learning, inspired from LSTM.
- RMLC network is better at extracting modal-specific information and long range cross-modal information, which stores modal information with memory cells. Meanwhile, RMLC sequentially exploits label correlations with the collaboration of both modal-specific information and cross-modal interaction.
- Extensive experiments on 5 benchmark multi-modal multi-label datasets verify the effectiveness of RMLC compared with several state-of-the-art approaches.

The rest of the paper is organized as follows. Section II briefly reviews some related works. Section III presents technical details of the proposed approach. Section IV reports experimental results, followed by the conclusion in Section V.

Corresponding author

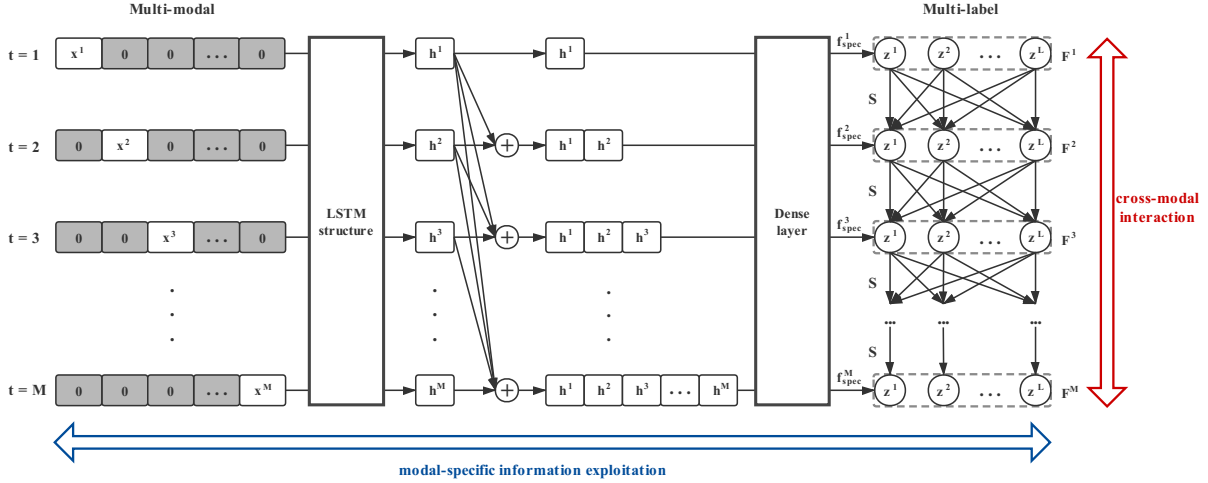


Fig. 1. The overall flowchart of RMLC approach, which is composed of a LSTM structure and a dense layer.  $\oplus$  denotes the concatenation the vectors. At  $t$ -th step, RMLC inputs modified feature vector to the LSTM structure. After obtaining output vector  $h^t$ , RMLC stacks all previous output vectors, i.e.,  $[h^1, h^2, \dots, h^t]$ . Based on the stacked output features of LSTM, RMLC learns label prediction from previous  $t$  modalities, which is denoted as  $f_{spec}^t$ . The LSTM layer is used for rethink process, which goes through  $M$  modalities sequentially and each step shares the same label correlation matrix  $\mathcal{S}$ . What's more, the label prediction at each step is propagated to the next step in the label prediction layer. RMLC uses the label prediction from previous step to polish the modal-specific label prediction, i.e., at  $t$ -th step, RMLC makes label prediction with collaboration between modal-specific label prediction  $f_{spec}^t$  and previous label predictions  $F^t$  of other labels.

## II. RELATED WORK

Our work is related to two branches of studies: multi-label learning and multi-modal learning. In this section, we briefly review some state-of-the-art approaches in the two fields.

Multi-label learning [3] [4] deals with objects annotated with multiple interdependent labels. Binary Relevance (BR) [5] is the most straightforward solution that decompose the multi-label problem into a set of independent binary classification tasks, but it neglects correlations among labels. To tackle label dependence, Classifier Chains (CC) [6] was proposed as a high-order approach to consider correlations between labels. However, the performance of CC is seriously affected by the training order of labels. So far, many approaches have been developed to improve the performance of multi-label learning by exploring various types of label correlations. For example, [7] exploits label correlations locally, [8] learns label-specific data representation for each label, [9] exploits global and local label correlations. However, the label correlations are simply obtained by common similarity measures, which may not be able to reflect complex relationships among labels. To address the above limitation, CAMEL [10] is proposed to learn label correlations via sparse reconstruction in the label space, which is capable of reflecting the collaborative relationships among labels regarding the final predictions.

Multi-modal learning aims to jointly utilize different information collected with diverse data collection techniques, e.g., [11] applied deep network to learn features over multiple modal data. Most approaches are mainly derived from the Canonical Correlation Analysis (CCA) methods, such as deep neural networks based CCA [12]. Meanwhile, Multi-modal multi-label learning has been widely studied, e.g.,

[9] exploits the consensus among different modalities, where multi-view latent spaces are correlated by HilbertSchmidt Independence Criterion (HSIC). It has been recognized that exploring individuality and commonality of heterogeneous features can further boost the performance of multi-modal data mining. SMISFL [13] jointly learns multiple modal-individual transformations and one sharable transformation; SIMM [14] leverages shared subspace exploitation and modal-specific information extraction; ICM2L [13] adopts an ensemble strategy to explicitly explore the individuality and commonality information of multi-label multiple view data in a unified model. Nevertheless, the above approaches rarely consider the label correlations. MMP [15] handles the consistencies among different views by requiring them to generate the same annotation result, and captures the correlations among different labels by imposing the similarity constraints. CS3G approach [16] handles types of interactions between multiple labels, while no interaction between features from different modalities. [17] 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 make each modality interacts and further reduce modal extraction cost, MCC [18] extends Classifier Chains to exploit label correlation with partial modalities.

## III. PROPOSED METHODOLOGY

In this section, we first summarize some formal notations used throughout this paper. And then we present detailedly the proposed method named Rethinking Modal-oriented Label Correlations (RMLC), along with its implementation. Fig. 1 shows the overall architecture of RMLC.

## A. Preliminaries

In the multi-modal multi-label learning, an instance is characterized with multiple modalities and annotated with multiple labels. Formally, let  $\mathcal{X} = \mathbb{R}^{d_1} \times \mathbb{R}^{d_2} \times \mathbb{R}^{d_M}$ , where  $d_m (1 \leq m \leq M)$  is the dimensionality of the  $m$ -th modality.

Suppose  $\mathcal{D} = \{(\mathbf{X}_i, \mathbf{Y}_i)\}_{i=1}^N$  represents a dataset with  $N$  samples. For the  $i$ -th instance,  $\mathbf{X}_i = [\mathbf{x}_i^1, \mathbf{x}_i^2, \dots, \mathbf{x}_i^M] \in \mathcal{X}$  is the feature vector with  $M$  modalities, where  $\mathbf{x}_i^m \in \mathbb{R}^{d_m}$ .  $\mathbf{Y}_i = [y_i^1, y_i^2, \dots, y_i^L] \in \mathcal{Y}$  is the corresponding label vector with  $L$  labels, where  $y_i^k \in \{-1, 1\}$ .  $y_i^k = 1 (k = 1, \dots, L)$  means the  $k$ -th label of  $\mathbf{X}_i$  is relevant;  $y_i^k = -1$  otherwise.

The task of multi-modal multi-label learning is to learn a predictive model:  $h : \mathcal{X} \rightarrow 2^{\mathcal{Y}}$  from  $\mathcal{D}$ , which can assign a set of proper labels for the unseen instance.

RMLC sequentially polishing the label prediction in a rethink manner and exploits label correlations based on the following two types of information: modal-specific information and cross-modal interaction. What's more, the final prediction of each label is composed of its own prediction and the prediction of other labels.

---

**Algorithm 1** The pseudo code of RMLC approach

---

**Input:**

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

**Output:**

- $F^M$ : multi-modal classifier trained with  $M$ -modalities

- 1: **repeat**
  - 2:   Initialize label correlation matrix  $\mathbf{S}$
  - 3:   Randomly select  $N_b$  instances from  $\mathcal{D}$
  - 4:   **for**  $t = 1 : M$  **do**
  - 5:     **for**  $i = 1 : N_b$  **do**
  - 6:       Modify  $\mathbf{X}_i$  to  $\hat{\mathbf{X}}_i^t$  with Eq. 1
  - 7:       Input  $\hat{\mathbf{X}}_i^t$  to the LSTM structure
  - 8:       Stack hidden output  $\mathbf{H}_i^t = [\mathbf{h}_i^0, \mathbf{h}_i^1, \dots, \mathbf{h}_i^t]$
  - 9:       Compute label prediction  $\mathbf{F}^t(\mathbf{H}_i^t)$  with Eq. 4
  - 10:       Calculate loss function  $\mathcal{L}_i^t$  with Eq. 5
  - 11:     **end for**
  - 12:     Calculate overall loss function  $\mathcal{L}^t$  with Eq. 6
  - 13:     Weight propagation: Obtain the derivative  $\partial \mathcal{L}^t / \partial \Phi$ ,  $\partial \mathcal{L}^t / \partial \mathbf{U}^t$ ,  $\partial \mathcal{L}^t / \partial \mathbf{b}^t$  and  $\partial \mathcal{L}^t / \partial \mathbf{S}$
  - 14:     Update parameters:  $\Phi$ ,  $\mathbf{U}^t$ ,  $\mathbf{b}^t$  and  $\mathbf{S}$
  - 15:   **end for**
  - 16: **until** converge
  - 17: **return**  $F^M$
- 

## B. Modal-specific information exploitation

It is well-known that each modality contains its own specific contribution to the multi-label prediction. Take image annotation as an example, a picture of pink rose may be represented with two modalities: HSV and Gist, while tagged with two labels *pink*, *flower* simultaneously [14]. Intuitively, we can infer *pink* from HSV (color) modality and *flower* from Gist (texture) modality. Existing approaches mainly try

to find the shared information between different modalities, while it is more reasonable to consider extracting their own specific information. At the  $t$ -th step, RMLC exploits specific information of the  $t$ -th individual modality  $\mathbf{x}_i^t$ . It is notable that the dimensionality of various modalities is heterogeneous, which is difficult to input to the LSTM structure.

First of all, we adapt  $\mathbf{X}_i$  as  $\hat{\mathbf{X}}_i^t = [\hat{\mathbf{x}}_i^1, \hat{\mathbf{x}}_i^2, \dots, \hat{\mathbf{x}}_i^M]$  according to Eq. 1.

$$\hat{\mathbf{x}}_i^m = \begin{cases} \mathbf{x}_i^m & m = t \\ \mathbf{0}^m & m \neq t \end{cases} \quad (1)$$

where  $\mathbf{0}^m \in \mathbb{R}^{d_m}$ . For example, if  $t = 4$ ,  $\hat{\mathbf{X}}_i^4 = [\mathbf{0}^1, \mathbf{0}^2, \mathbf{0}^3, \mathbf{x}_i^4, \dots, \mathbf{0}^M]$ .

Secondly, we input the modified  $\hat{\mathbf{X}}_i^t$  to LSTM structure, which includes input gate, forget gate, cell state, output gate, and output vector.

$$\begin{aligned} \mathbf{i}_i^t &= \sigma(\mathbf{W}_i[\hat{\mathbf{X}}_i^t, \mathbf{h}_i^{t-1}, \mathbf{c}_i^{t-1}] + \mathbf{b}_i) \\ \mathbf{f}_i^t &= \sigma(\mathbf{W}_f[\hat{\mathbf{X}}_i^t, \mathbf{h}_i^{t-1}, \mathbf{c}_i^{t-1}] + \mathbf{b}_f) \\ \mathbf{c}_i^t &= \mathbf{f}_i^t \mathbf{c}_i^{t-1} + \mathbf{i}_i^t \tanh(\mathbf{W}_c[\hat{\mathbf{X}}_i^t, \mathbf{h}_i^{t-1}] + \mathbf{b}_c) \\ \mathbf{o}_i^t &= \sigma(\mathbf{W}_o[\hat{\mathbf{X}}_i^t, \mathbf{h}_i^{t-1}, \mathbf{c}_i^t] + \mathbf{b}_o) \\ \mathbf{h}_i^t &= \mathbf{o}_i^t \tanh(\mathbf{c}_i^t) \end{aligned} \quad (2)$$

where  $\mathbf{W}_i, \mathbf{W}_f, \mathbf{W}_c, \mathbf{W}_o \in \mathbb{R}^h$  represents the weight matrices from the cell to gate vectors,  $\mathbf{b}_i, \mathbf{b}_f, \mathbf{b}_c, \mathbf{b}_o \in \mathbb{R}^h$  denotes bias vectors. For simplicity, we denote all the parameters in the LSTM structure as  $\Phi$ .

In order to combine the specific information extracted from previous  $t$  modalities  $\{\mathbf{x}_i^1, \mathbf{x}_i^2, \dots, \mathbf{x}_i^t\}$ , we stack all the output features, which is denoted as  $\mathbf{H}_i^t = [\mathbf{h}_i^1, \mathbf{h}_i^2, \dots, \mathbf{h}_i^t]$  and  $\mathbf{H}_i^t \in \mathbb{R}^{1 \times (t \cdot h)}$ .

Last but not the least, we add a dense layer between stacked output features and label prediction layer. The dense layer learns multi-label embedding to transform the output features  $\mathbf{H}_i^t$  to label vector, which can be predicted by Eq. 3.

$$\begin{aligned} \mathbf{f}_{spec}^t(\mathbf{H}_i^t) &= \mathbf{U}^t \mathbf{H}_i^t + \mathbf{b}^t \\ &= \mathbf{U}^t [\mathbf{h}_i^1, \mathbf{h}_i^2, \dots, \mathbf{h}_i^t] + \mathbf{b}^t \end{aligned} \quad (3)$$

where  $\mathbf{U}^t \in \mathbb{R}^{(t \cdot h) \times L}$  is the weight vector,  $\mathbf{b}^t$  is the bias vector, and  $\mathbf{f}_{spec}^t(\cdot) \in \mathbb{R}^{1 \times L}$ .

## C. Cross-modal interaction

The exploitation of cross-modal interaction can be considered as a sequential case. To characterize the collaborative relationship among multiple labels regarding the final prediction, RMLC also learns label correlation matrix  $\mathbf{S} = [s_{jk}]_{L \times L}$ , where  $s_{jk}$  reflects the contribution of the  $j$ -th label with all previous  $t - 1$  modalities to the  $k$ -th label.

Each step in LSTM structure represents a rethink step, which classifies labels and exploit label correlations simultaneously. We denote the predicted label vector as  $\mathbf{Z}_i^t = [z_i^{1,t}, z_i^{2,t}, \dots, z_i^{L,t}]$  and we initialize  $\mathbf{Z}_i^0 = \mathbf{0}$  for all instances, where  $\mathbf{Z}_i^t = \mathbf{F}_i^t(\mathbf{H}_i^t)$  according to Eq. 4. With the help of

TABLE I

EXAMPLE OF LABEL PREDICTOR  $\mathbf{f}_{spec}^t(\cdot)$  AND  $\mathbf{F}^t(\cdot)$ , WHERE  $\mathbf{f}_{spec}^t(\cdot)$  DENOTES LABEL PREDICTION WITH THE MODAL-SPECIFIC INFORMATION EXPLOITATION,  $\mathbf{F}^t(\cdot)$  DENOTES FINAL LABEL PREDICTION WHICH INVOLVES COLLABORATION BETWEEN ITS OWN PREDICTION AND THE PREDICTION OF OTHER LABELS FROM THE LAST STEP.  $\mathbf{S}$  DENOTES LABEL CORRELATION MATRIX.

step	label predictor with modal-specific information exploitation	label predictor with both modal-specific and cross-modal interaction
$t = 1$	$\mathbf{f}_{spec}^1(\mathbf{H}_i^1) = \mathbf{U}^1 \mathbf{h}_i^1 + \mathbf{b}^1$	$\mathbf{F}^1(\mathbf{H}_i^1) = \sigma(\mathbf{U}^1 \mathbf{h}_i^1 + \mathbf{b}^1)$
$t = 2$	$\mathbf{f}_{spec}^2(\mathbf{H}_i^2) = \mathbf{U}^2 [\mathbf{h}_i^1, \mathbf{h}_i^2] + \mathbf{b}^2$	$\mathbf{F}^2(\mathbf{H}_i^2) = \sigma(\mathbf{U}^2 [\mathbf{h}_i^1, \mathbf{h}_i^2] + \mathbf{b}^2 + \mathbf{F}^1(\mathbf{h}_i^1) \mathbf{S})$
$t = 3$	$\mathbf{f}_{spec}^3(\mathbf{H}_i^3) = \mathbf{U}^3 [\mathbf{h}_i^1, \mathbf{h}_i^2, \mathbf{h}_i^3] + \mathbf{b}^3$	$\mathbf{F}^3(\mathbf{H}_i^3) = \sigma(\mathbf{U}^3 [\mathbf{h}_i^1, \mathbf{h}_i^2, \mathbf{h}_i^3] + \mathbf{b}^3 + \mathbf{F}^2([\mathbf{h}_i^1, \mathbf{h}_i^2]) \mathbf{S})$
...	...	...
$t = M$	$\mathbf{f}_{spec}^M(\mathbf{H}_i^M) = \mathbf{U}^M [\mathbf{h}_i^1, \mathbf{h}_i^2, \dots, \mathbf{h}_i^M] + \mathbf{b}^M$	$\mathbf{F}^M(\mathbf{H}_i^M) = \sigma(\mathbf{U}^M [\mathbf{h}_i^1, \mathbf{h}_i^2, \dots, \mathbf{h}_i^M] + \mathbf{b}^M + \mathbf{F}^{M-1}([\mathbf{h}_i^1, \mathbf{h}_i^2, \dots, \mathbf{h}_i^{M-1}]) \mathbf{S})$
Output	$\mathbf{f}_{spec}^M(\cdot)$	$\mathbf{F}^M(\cdot)$

TABLE II

CHARACTERISTIC OF THE BENCHMARK MULTI-MODAL MULTI-LABEL DATASETS.  $N$ ,  $M$  AND  $L$  DENOTE THE NUMBER OF INSTANCES, MODALITIES AND LABELS IN EACH DATASET, RESPECTIVELY.  $d_m$  SHOWS THE DIMENSIONALITY OF EACH MODALITY.

dataset	# of size $N$	# of modalities $M$	# of labels $L$	# the dimensionality for each modality $d_m$
<i>ML2000</i>	2000	3	5	[500,1040,576]
<i>MSRC</i>	591	3	24	[500,1040,576]
<i>Taobao</i>	2079	4	30	[500,48,81,24]
<i>FCVID</i>	4388	5	28	[400,400,400,400,400]
<i>MSRA</i>	15000	7	50	[256,225,64,144,75,128,7]

memory in the LSTM structure, each  $\mathbf{Z}_i^t$  is passed down to  $(t + 1)$ -th step.

Combine both modal-specific prediction and cross-modal interactions together, we predict labels according to a nonlinear function  $\mathbf{F}^t(\cdot)$ .

$$\mathbf{F}^t(\mathbf{H}_i^t) = \sigma(\mathbf{f}_{spec}(\mathbf{H}_i^t) + \mathbf{Z}_i^{t-1} \mathbf{S}) \quad (4)$$

where  $\sigma(\cdot)$  is the sigmoid function,  $\mathbf{Z}_i^{t-1} = \mathbf{F}^{t-1}(\mathbf{H}_i^{t-1})$  denotes the predicted label vector at  $t - 1$ -th step. The modal-specific term  $\mathbf{f}_{spec}^t(\cdot)$  concentrates on extracting individual information of each modality to predict labels independently, which is similar to BR. The memory term  $\mathbf{Z}_i^{t-1} \mathbf{S}$  transforms the previous prediction to the current label vector space, which can also be considered as the exploitation for label correlations with cross-modal interactions. Table I shows example of label predictor  $\mathbf{f}_{spec}^t(\cdot)$  and  $\mathbf{F}^t(\cdot)$  at each step.

#### D. Loss function

As a general training procedure, it focuses on reducing the errors made in the current status of the network. Thus we design weight binary cross-entropy loss function for label prediction at the  $t$ -th step according to:

$$\mathcal{L}_i^t = - \sum_{k=1}^L (y_i^k \log z_i^{k,t} + (1 - y_i^k) \log(1 - z_i^{k,t})) \quad (5)$$

where  $\mathbf{Z}_i^t = [z_i^{1,t}, z_i^{2,t}, \dots, z_i^{L,t}]$  is predicted by  $\mathbf{F}^t(\mathbf{H}_i^t)$  and  $z_i^{k,t}$  is the prediction of the  $k$ -th label at the  $t$ -th step.

Above all, we combine label loss function  $\mathcal{L}_i^t$  and regularization term  $\|\mathbf{U}^t\|_2^2$  to calculate the overall loss function:

$$\mathcal{L}^t = \sum_i^{N_b} \mathcal{L}_i^t + \lambda \|\mathbf{U}^t\|_2^2 \quad (6)$$

where  $\|\cdot\|_2$  represents  $L_2$  norm and  $\lambda$  is the trade-off between the label loss function and the regularization term.

In the training procedure, the derivatives are taken with the help of back propagation technique, and the pseudo code of RMLC is summarized in Algorithm 1. At  $t$ -th step, we adopt the popular optimization algorithm Adam [19] to update all the parameters in  $\Phi$ ,  $\mathbf{U}^t$ ,  $\mathbf{b}^t$  and  $\mathbf{S}$  simultaneously.

## IV. EXPERIMENTS

In this section, we conduct extensive experiments on various datasets to validate the effectiveness of RMLC.

### A. Dataset description

For comprehensive performance evaluation, 5 benchmark multi-modal multi-label datasets are collected as follows. Table II summarizes the detailed characteristics of these datasets, which are organized in ascending order of  $M$  (the number of modalities).

- *ML2000* [20] is an image dataset with 2000 images from 5 categories (desert, mountains, sea, sunset and trees), and we extract 3 types of modalities for each image including: BoW, FV and HOG.
- *MSRC* [21] is used for object class recognition. Similar to *ML2000*, we extract 3 modalities of features: BoW, FV and HOG for each image. Each image may be affiliated to two or more labels. Given different modalities of description for an image, the task is to predict its all possible labels.
- *Taobao* [16] is used for shopping items classification which has 2079 instances and 30 labels. Description images of items are crawled from a shopping website, and four types of features, e.g., BoW, Gabor, HOG,



TABLE III  
COMPARISON RESULTS (MEAN  $\pm$  STANDARD DEVIATION) OF RMLC 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.

Datasets	Approaches	Evaluation Metrics					
		Hamming Loss $\downarrow$	Ranking Loss $\downarrow$	Subset Accuracy $\uparrow$	Macro F1 $\uparrow$	Example F1 $\uparrow$	Micro F1 $\uparrow$
ML2000	CAMEL(B)	0.089 $\pm$ 0.011	0.063 $\pm$ 0.011	0.655 $\pm$ 0.040	0.804 $\pm$ 0.025	0.769 $\pm$ 0.031	0.806 $\pm$ 0.024
	CAMEL(C)	0.098 $\pm$ 0.011	0.067 $\pm$ 0.010	0.633 $\pm$ 0.036	0.781 $\pm$ 0.023	0.739 $\pm$ 0.029	0.783 $\pm$ 0.024
	DMP	0.103 $\pm$ 0.010	0.074 $\pm$ 0.010	0.617 $\pm$ 0.031	0.781 $\pm$ 0.020	0.753 $\pm$ 0.026	0.783 $\pm$ 0.021
	CS3G	0.119 $\pm$ 0.010	0.092 $\pm$ 0.013	0.564 $\pm$ 0.033	0.738 $\pm$ 0.021	0.697 $\pm$ 0.029	0.744 $\pm$ 0.021
	MCC	0.105 $\pm$ 0.012	0.082 $\pm$ 0.011	0.662 $\pm$ 0.032	0.780 $\pm$ 0.027	0.787 $\pm$ 0.022	0.784 $\pm$ 0.024
	RMLC	<b>0.080<math>\pm</math>0.007</b>	<b>0.058<math>\pm</math>0.009</b>	<b>0.753<math>\pm</math>0.024</b>	<b>0.835<math>\pm</math>0.016</b>	<b>0.830<math>\pm</math>0.021</b>	<b>0.834<math>\pm</math>0.016</b>
MSRC	CAMEL(B)	0.059 $\pm$ 0.009	0.033 $\pm$ 0.009	0.322 $\pm$ 0.058	0.680 $\pm$ 0.043	0.805 $\pm$ 0.032	0.814 $\pm$ 0.029
	CAMEL(C)	0.106 $\pm$ 0.020	0.153 $\pm$ 0.074	0.070 $\pm$ 0.073	0.208 $\pm$ 0.041	0.618 $\pm$ 0.074	0.629 $\pm$ 0.071
	DMP	0.067 $\pm$ 0.008	0.047 $\pm$ 0.010	0.291 $\pm$ 0.033	0.643 $\pm$ 0.045	0.789 $\pm$ 0.028	0.796 $\pm$ 0.025
	CS3G	0.077 $\pm$ 0.008	0.043 $\pm$ 0.010	0.205 $\pm$ 0.043	0.506 $\pm$ 0.040	0.729 $\pm$ 0.024	0.744 $\pm$ 0.026
	MCC	0.068 $\pm$ 0.008	0.054 $\pm$ 0.011	0.345 $\pm$ 0.088	0.652 $\pm$ 0.024	0.798 $\pm$ 0.024	0.799 $\pm$ 0.023
	RMLC	<b>0.048<math>\pm</math>0.010</b>	<b>0.031<math>\pm</math>0.010</b>	<b>0.472<math>\pm</math>0.048</b>	<b>0.761<math>\pm</math>0.059</b>	<b>0.854<math>\pm</math>0.035</b>	<b>0.856<math>\pm</math>0.032</b>
Taobao	CAMEL(B)	0.032 $\pm$ 0.001	0.151 $\pm$ 0.014	0.150 $\pm$ 0.020	0.034 $\pm$ 0.008	0.054 $\pm$ 0.012	0.101 $\pm$ 0.023
	CAMEL(C)	0.033 $\pm$ 0.002	0.159 $\pm$ 0.012	0.238 $\pm$ 0.046	0.182 $\pm$ 0.028	0.266 $\pm$ 0.044	0.373 $\pm$ 0.053
	DMP	0.053 $\pm$ 0.002	0.238 $\pm$ 0.020	0.218 $\pm$ 0.017	0.231 $\pm$ 0.015	0.336 $\pm$ 0.014	0.359 $\pm$ 0.015
	CS3G	0.064 $\pm$ 0.003	0.171 $\pm$ 0.009	0.104 $\pm$ 0.017	0.125 $\pm$ 0.015	0.323 $\pm$ 0.025	0.332 $\pm$ 0.026
	MCC	0.054 $\pm$ 0.002	0.235 $\pm$ 0.016	0.213 $\pm$ 0.023	0.222 $\pm$ 0.026	0.333 $\pm$ 0.026	0.354 $\pm$ 0.021
	RMLC	<b>0.028<math>\pm</math>0.002</b>	<b>0.115<math>\pm</math>0.020</b>	<b>0.476<math>\pm</math>0.035</b>	<b>0.352<math>\pm</math>0.027</b>	<b>0.477<math>\pm</math>0.038</b>	<b>0.555<math>\pm</math>0.039</b>
FCVID	CAMEL(B)	0.018 $\pm$ 0.001	0.027 $\pm$ 0.006	0.534 $\pm$ 0.030	0.697 $\pm$ 0.022	0.551 $\pm$ 0.031	0.695 $\pm$ 0.024
	CAMEL(C)	0.020 $\pm$ 0.001	0.031 $\pm$ 0.005	0.485 $\pm$ 0.019	0.659 $\pm$ 0.017	0.503 $\pm$ 0.017	0.658 $\pm$ 0.016
	DMP	0.027 $\pm$ 0.002	0.052 $\pm$ 0.007	0.520 $\pm$ 0.033	0.660 $\pm$ 0.023	0.642 $\pm$ 0.026	0.658 $\pm$ 0.023
	CS3G	0.020 $\pm$ 0.001	0.027 $\pm$ 0.003	0.571 $\pm$ 0.018	0.693 $\pm$ 0.014	0.673 $\pm$ 0.022	0.720 $\pm$ 0.015
	MCC	0.027 $\pm$ 0.001	0.052 $\pm$ 0.005	0.522 $\pm$ 0.024	0.667 $\pm$ 0.009	0.647 $\pm$ 0.019	0.663 $\pm$ 0.014
	RMLC	<b>0.013<math>\pm</math>0.001</b>	<b>0.023<math>\pm</math>0.003</b>	<b>0.757<math>\pm</math>0.023</b>	<b>0.829<math>\pm</math>0.015</b>	<b>0.797<math>\pm</math>0.023</b>	<b>0.821<math>\pm</math>0.017</b>
MSRA	CAMEL(B)	0.046 $\pm$ 0.001	0.159 $\pm$ 0.010	0.057 $\pm$ 0.010	0.069 $\pm$ 0.004	0.232 $\pm$ 0.010	0.329 $\pm$ 0.014
	CAMEL(C)	<b>0.045<math>\pm</math>0.001</b>	0.154 $\pm$ 0.009	0.066 $\pm$ 0.010	0.079 $\pm$ 0.002	0.248 $\pm$ 0.006	0.349 $\pm$ 0.010
	DMP	0.046 $\pm$ 0.001	0.204 $\pm$ 0.005	0.053 $\pm$ 0.005	0.054 $\pm$ 0.002	0.216 $\pm$ 0.007	0.311 $\pm$ 0.009
	CS3G	0.050 $\pm$ 0.001	0.140 $\pm$ 0.007	0.061 $\pm$ 0.009	0.041 $\pm$ 0.007	0.273 $\pm$ 0.011	0.324 $\pm$ 0.018
	MCC	0.048 $\pm$ 0.001	0.195 $\pm$ 0.005	0.076 $\pm$ 0.003	0.073 $\pm$ 0.004	0.266 $\pm$ 0.005	0.359 $\pm$ 0.006
	RMLC	0.048 $\pm$ 0.001	<b>0.132<math>\pm</math>0.005</b>	<b>0.136<math>\pm</math>0.004</b>	<b>0.186<math>\pm</math>0.007</b>	<b>0.334<math>\pm</math>0.011</b>	<b>0.421<math>\pm</math>0.009</b>

HSVHist, are extracted to construct 4 modalities of data. Corresponding categories path of an item provides the label sets.

- *FCVID* [22] 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. Given different modalities of description for a video, the task is to predict its possible categories.
- *MSRA* is a subset of a salient object recognition dataset [23], 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.

### B. Evaluation metrics

For performance evaluation, we use 6 widely-adopted evaluation metrics, including *Hamming Loss*, *Ranking Loss*, *Subset Accuracy*, *Macro F1*, *Example F1* and *Micro F1* [3], which

consider the performance of multi-label predictor from various aspects. All the employed approaches vary within the interval  $[0, 1]$ . For the first two evaluation metrics, the smaller values indicate the better performance, while for the last four evaluation metrics, the larger indicate the better performance.

### C. Compared approaches

The performance of RMLC is compared against 5 multi-modal multi-label learning approaches, listed as follows:

- CAMEL(B) & CAMEL(C): CAMEL [10] is a state-of-the-art multi-label approach, which learn the label correlations via sparse reconstruction in the label space. CAMEL(B) stands for the best performance obtained from the best single modality. CAMEL(C) stands for concatenating all modalities as a single modal input.
- DMP [24]: A multi-modal approach which automatically extracts instance-specifically discriminative modal sequence for reducing the cost of feature extraction. Here, we treat each label independently, i.e., for each label, DMP trains classifiers using different modalities.
- CS3G [16]: A multi-modal multi-label approach utilizing multi-modal information in a privacy-preserving style to

TABLE IV

COMPARISON RESULTS (MEAN  $\pm$  STANDARD DEVIATION) OF RMLC\_NS AND RMLC, WHERE RMLC\_NS DENOTES THE MODEL WITHOUT CONSIDERING LABEL CORRELATION MATRIX  $\mathbf{S}$ . 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 $\uparrow$
ML2000	RMLC_NS	0.083 $\pm$ 0.008	0.059 $\pm$ 0.009	0.742 $\pm$ 0.028	0.828 $\pm$ 0.017	0.821 $\pm$ 0.019	0.826 $\pm$ 0.017
	RMLC	<b>0.080<math>\pm</math>0.007</b>	<b>0.058<math>\pm</math>0.009</b>	<b>0.753<math>\pm</math>0.024</b>	<b>0.835<math>\pm</math>0.016</b>	<b>0.830<math>\pm</math>0.021</b>	<b>0.834<math>\pm</math>0.016</b>
MSRC	RMLC_NS	0.059 $\pm$ 0.006	0.040 $\pm$ 0.012	0.359 $\pm$ 0.062	0.693 $\pm$ 0.058	0.815 $\pm$ 0.022	0.819 $\pm$ 0.022
	RMLC	<b>0.048<math>\pm</math>0.010</b>	<b>0.031<math>\pm</math>0.010</b>	<b>0.472<math>\pm</math>0.048</b>	<b>0.761<math>\pm</math>0.059</b>	<b>0.854<math>\pm</math>0.035</b>	<b>0.856<math>\pm</math>0.032</b>
Taobao	RMLC_NS	0.030 $\pm$ 0.002	0.137 $\pm$ 0.017	0.394 $\pm$ 0.039	0.284 $\pm$ 0.027	0.383 $\pm$ 0.038	0.486 $\pm$ 0.041
	RMLC	<b>0.028<math>\pm</math>0.002</b>	<b>0.115<math>\pm</math>0.020</b>	<b>0.476<math>\pm</math>0.035</b>	<b>0.352<math>\pm</math>0.027</b>	<b>0.477<math>\pm</math>0.038</b>	<b>0.555<math>\pm</math>0.039</b>
FCVID	RMLC_NS	0.013 $\pm$ 0.001	0.025 $\pm$ 0.004	0.702 $\pm$ 0.021	0.804 $\pm$ 0.016	0.703 $\pm$ 0.019	0.801 $\pm$ 0.017
	RMLC	<b>0.013<math>\pm</math>0.001</b>	<b>0.023<math>\pm</math>0.003</b>	<b>0.757<math>\pm</math>0.023</b>	<b>0.829<math>\pm</math>0.015</b>	<b>0.797<math>\pm</math>0.023</b>	<b>0.821<math>\pm</math>0.017</b>
MSRA	RMLC_NS	<b>0.045<math>\pm</math>0.000</b>	<b>0.115<math>\pm</math>0.005</b>	0.119 $\pm$ 0.007	0.125 $\pm$ 0.009	0.280 $\pm$ 0.008	0.393 $\pm$ 0.009
	RMLC	0.048 $\pm$ 0.001	0.132 $\pm$ 0.005	<b>0.136<math>\pm</math>0.004</b>	<b>0.186<math>\pm</math>0.007</b>	<b>0.334<math>\pm</math>0.011</b>	<b>0.421<math>\pm</math>0.009</b>

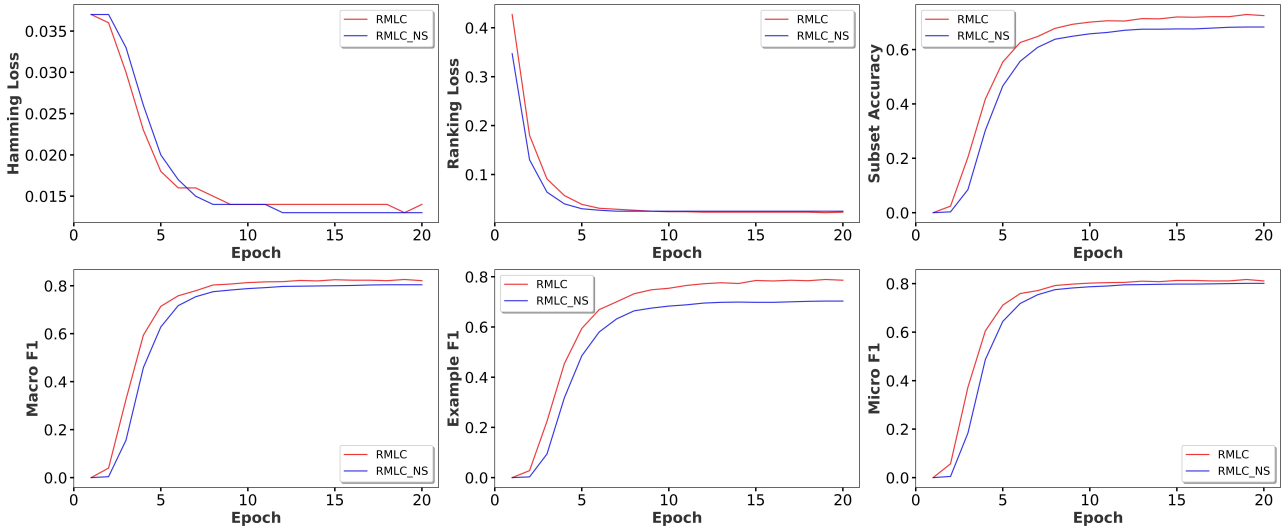


Fig. 2. Convergence analysis of RMLC and RMLC\_NS approach on the FCVID dataset.

deal with multi-label tasks. CS3G treats each modality unequally and has the ability of extracting the most useful modal features for final recommendation as well.

- MCC [18]: A novel multi-modal multi-label approach considering not only interrelation among different modalities, but also relationship among different labels. And MCC makes convince prediction with partial modalities.

#### D. Experimental results

For each dataset, we perform 10-fold cross-validation and take the mean metrics results, standard deviations for all compared approaches. Detailed experimental results are reported in Table III, which obviously shows RMLC outperforms the other 5 compared approaches on all evaluation metrics. We set trade-off parameter  $\lambda = 0.1$  and batch size  $N_b = 64$ .

Based on the experimental results in Table III, we obtain the following observations:

- From the results of CAMLE(C) approach, it is obvious shown that roughly concatenating all modalities as a

single modality may not always be the wise choice, which indicates that it is necessary to extract information of each individual modality separately.

- RMLC approach achieves the best performance compared with several state-of-the-art approaches on the benchmark datasets, which shows priority of our proposed RMLC approach in multi-modal multi-label learning problem. The main reason is that RMLC takes advantage of multiple modal-specific information, as well as exploiting label correlations to polish multi-label prediction with the help of cross-modal interaction.

1) *Influence of label correlations*: RMLC improves multi-label prediction in a rethink manner, which makes use of previous modal information to better characterize the relationship among different labels. In order to validate the effectiveness of label correlations across different modalities, we keep the basic structure of RMLC and only adopt  $\sigma(\mathbf{f}_{spec}^M)$  as the final prediction, which is denoted as RMLC\_NS. As shown in Table IV, RMLC performs better than RMLC\_NS, which

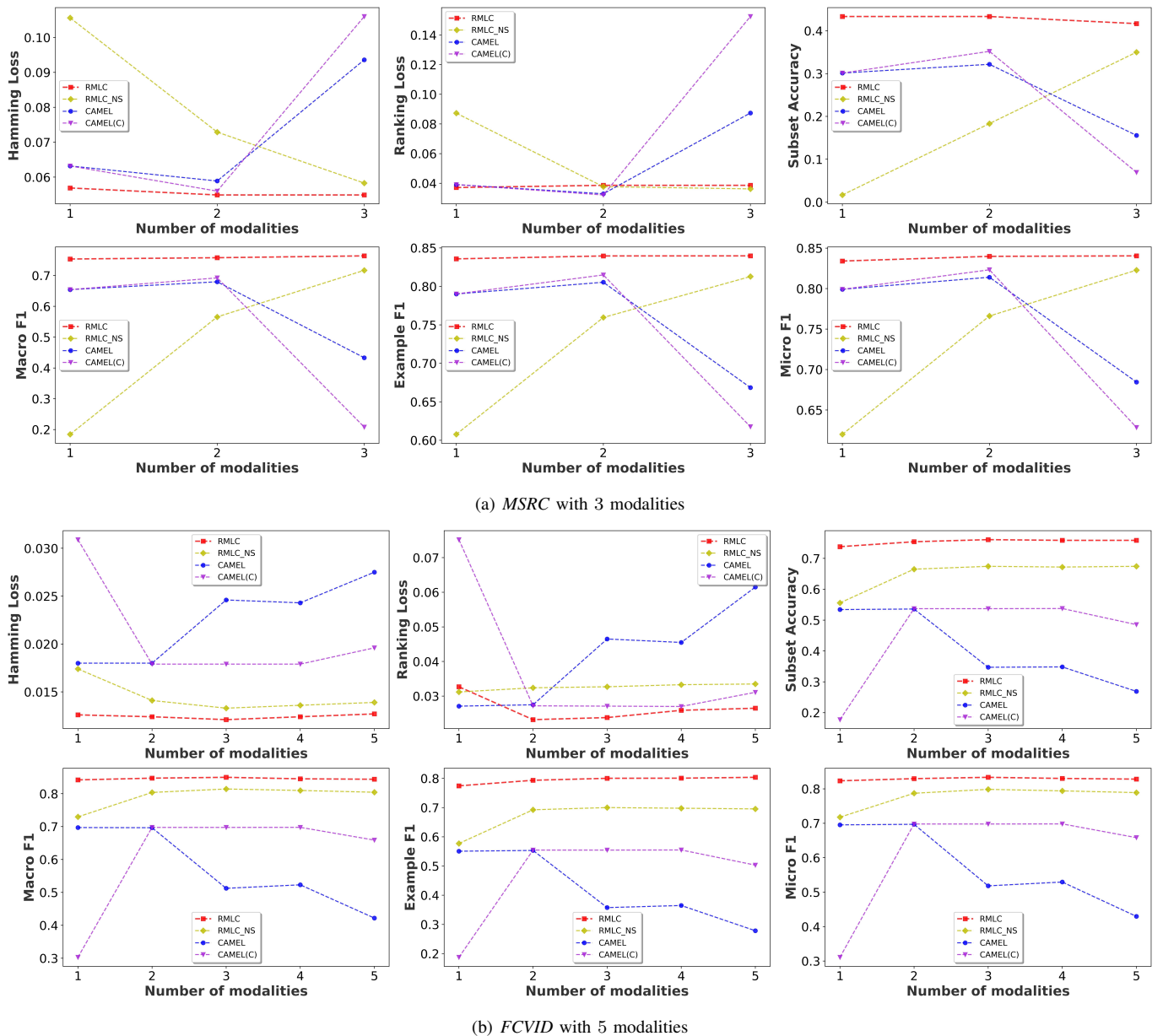


Fig. 3. Performance of RMLC, RMLC\_NS, CAMEL and CAMEL(C) with increase of the modality on the MSRC and FCVID dataset. At  $t$ -th step, CAMEL adopts  $t$ -th modality as input, while CAMEL(C) concatenates  $t$  previous modalities as input.

demonstrates the effectiveness of exploiting label correlations. In other words, it is not enough to merely fuse modal-specific information. RMLC is of great effectiveness to integrate the learned label correlations into the desired label prediction.

Nevertheless, RMLC\_NS performs better than other state-of-the-art approaches shown in Table III, which validates the effectiveness of extracting specific information of each individual modalities sequentially.

2) *Convergence*: We conduct experiments to investigate convergence of RMLC. Due to page limit, we only report the experimental results on the FCVID datasets on all evaluation metrics. The curves in Fig. 2 illustrate the change of metrics results with the number of epoch. It is obvious that both RMLC

and RMLC\_NS can converge quickly with a few number of epochs, while RMLC performs better than RMLC\_NS.

3) *Influence of modal-specific information exploitation*: To further examine the effectiveness of modal-specific information exploitation, we output label prediction of RMLC at each step. In addition, we conduct experiment on each modality with state-of-the-art multi-label approach CAMEL. In Fig. 3, the curves show the change of the metric results of RMLC, RMLC\_NS, CAMEL and CAMEL(C) with previous  $t$  modalities. The experimental results in Fig. 3 can not only validate the favorable performance of RMLC in exploiting each individual modality, but also further confirm the contribution of previous  $t$  modalities.

## V. CONCLUSION

In this paper, a novel end-to-end network based approach is proposed to solve the multi-modal multi-label problem. Specifically, we enhance the communication among different modalities while remaining modal-specific characteristics. What's more, we exploit label correlations with the collaboration of modal-specific and cross-modal interaction. Experiments on several benchmark multi-modal multi-label datasets demonstrate the superiority of our proposed RMLC over related competitive approaches. Different modalities have different degree of consistency and complementarity, thus it will be an interesting work to exploit label correlations into reinforcement learning environment in the future.

## ACKNOWLEDGMENT

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] Q. Tan, G. Yu, J. Wang, C. Domeniconi, and X. Zhang, "Individuality and commonality-based multiview multilabel learning," *IEEE transactions on cybernetics*, 2019.
- [2] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [3] M.-L. Zhang and Z.-H. Zhou, "A review on multi-label learning algorithms," *IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING*, vol. 26, no. 8, p. 1819, 2014.
- [4] E. Gibaja and S. Ventura, "A tutorial on multilabel learning," *ACM Computing Surveys (CSUR)*, vol. 47, no. 3, p. 52, 2015.
- [5] M. R. Boutell, J. Luo, X. Shen, and C. M. Brown, "Learning multi-label scene classification," *Pattern recognition*, vol. 37, no. 9, pp. 1757–1771, 2004.
- [6] J. Read, B. Pfahringer, G. Holmes, and E. Frank, "Classifier chains for multi-label classification," *Machine learning*, vol. 85, no. 3, p. 333, 2011.
- [7] S.-J. Huang and Z.-H. Zhou, "Multi-label learning by exploiting label correlations locally," in *Twenty-sixth AAAI conference on artificial intelligence*, 2012.
- [8] J. Huang, G. Li, Q. Huang, and X. Wu, "Learning label-specific features and class-dependent labels for multi-label classification," *IEEE transactions on knowledge and data engineering*, vol. 28, no. 12, pp. 3309–3323, 2016.
- [9] Y. Zhu, J. T. Kwok, and Z.-H. Zhou, "Multi-label learning with global and local label correlation," *IEEE Transactions on Knowledge & Data Engineering*, no. 6, pp. 1081–1094, 2018.
- [10] L. Feng, B. An, and S. He, "Collaboration based multi-label learning," in *Thirty-Third AAAI Conference on Artificial Intelligence*, 2019, pp. 3550–3557.
- [11] R. Kiros, R. Salakhutdinov, and R. Zemel, "Multimodal neural language models," in *International Conference on Machine Learning*, 2014, pp. 595–603.
- [12] W. Wang, R. Arora, K. Livescu, and J. Bilmes, "On deep multi-view representation learning," in *International Conference on Machine Learning*, 2015, pp. 1083–1092.
- [13] F. Wu, X.-Y. Jing, J. Zhou, Y. Ji, C. Lan, Q. Huang, and R. Wang, "Semi-supervised multi-view individual and sharable feature learning for webpage classification," in *The World Wide Web Conference*. ACM, 2019, pp. 3349–3355.
- [14] X. Wu, Q.-G. Chen, Y. Hu, D. Wang, X. Chang, X. Wang, and M.-L. Zhang, "Multi-view multi-label learning with view-specific information extraction," in *Proceedings of the 28th International Joint Conference on Artificial Intelligence*. AAAI Press, 2019, pp. 3884–3890.
- [15] Z. He, C. Chen, J. Bu, P. Li, and D. Cai, "Multi-view based multi-label propagation for image annotation," vol. 168, no. C, pp. 853–860, 2015.
- [16] H.-J. Ye, D.-C. Zhan, X. Li, Z.-C. Huang, and Y. Jiang, "College student scholarships and subsidies granting: A multi-modal multi-label approach," in *2016 IEEE 16th International Conference on Data Mining (ICDM)*. IEEE, 2016, pp. 559–568.
- [17] Y. Xing, G. Yu, C. Domeniconi, J. Wang, and Z. Zhang, "Multi-label co-training," in *Proceedings of the 27th International Joint Conference on Artificial Intelligence*. AAAI Press, 2018, pp. 2882–2888.
- [18] Y. Zhang, C. Zeng, H. Cheng, C. Wang, and L. 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)*. IEEE, 2019, pp. 838–843.
- [19] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proceedings of the 3rd International Conference on Learning Representations*, 2015.
- [20] M. L. Zhang and Z. H. Zhou, "Ml-knn: A lazy learning approach to multi-label learning," *Pattern Recognition*, vol. 40, no. 7, pp. 2038–2048, 2007.
- [21] F. Schroff, A. Criminisi, and A. Zisserman, "Harvesting image databases from the web," *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 4, pp. 754–766, 2010.
- [22] Y.-G. Jiang, Z. Wu, J. Wang, X. Xue, and S.-F. Chang, "Exploiting feature and class relationships in video categorization with regularized deep neural networks," *IEEE Transactions on Pattern Analysis & Machine Intelligence*, vol. 40, no. 2, pp. 352–364, 2018.
- [23] T. Liu, Z. Yuan, J. Sun, J. Wang, N. Zheng, X. Tang, and H.-Y. Shum, "Learning to detect a salient object," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 2, pp. 353–367, 2010.
- [24] Y. Yang, D.-C. Zhan, Y. Fan, and Y. Jiang, "Instance specific discriminative modal pursuit: A serialized approach," in *Asian Conference on Machine Learning*, 2017, pp. 65–80.