



## Annotation Efficient Cross-Modal Retrieval with Adversarial Attentive Alignment

---

Po-Yao Huang, Guoliang Kang, Wenhe Liu, Xiaojun Chang and  
Alexander Hauptmann

EasyChair preprints are intended for rapid dissemination of research results and are integrated with the rest of EasyChair.

February 3, 2021

# Annotation Efficient Cross-Modal Retrieval with Adversarial Attentive Alignment

Po-Yao Huang<sup>1</sup>, Guoliang Kang<sup>1</sup>, Wenhe Liu<sup>1</sup>, Xiaojun Chang<sup>2\*</sup>, and Alexander G. Hauptmann<sup>1</sup>  
poyaoh,gkang,wenhel@cs.cmu.edu,cxj273@gmail.com,alex@cs.cmu.edu

<sup>1</sup>Language Technologies Institute, Carnegie Mellon University

<sup>2</sup>Faculty of Information Technology, Monash University

## ABSTRACT

Visual-semantic embeddings are central to many multimedia applications such as cross-modal retrieval between visual data and natural language descriptions. Conventionally, learning a joint embedding space relies on large parallel multimodal corpora. Since massive human annotation is expensive to obtain, there is a strong motivation in developing versatile algorithms to learn from large corpora with fewer annotations. In this paper, we propose a novel framework to leverage automatically extracted regional semantics from un-annotated images as additional weak supervision to learn visual-semantic embeddings. The proposed model employs adversarial attentive alignments to close the inherent heterogeneous gaps between annotated and un-annotated portions of visual and textual domains. To demonstrate its superiority, we conduct extensive experiments on sparsely annotated multimodal corpora. The experimental results show that the proposed model outperforms state-of-the-art visual-semantic embedding models by a significant margin for cross-modal retrieval tasks on the sparse Flickr30k and MS-COCO datasets. It is also worth noting that, despite using only 20% of the annotations, the proposed model can achieve competitive performance (Recall at 10 > 80.0% for 1K and > 70.0% for 5K text-to-image retrieval) compared to the benchmarks trained with the complete annotations.

## KEYWORDS

Cross-modal Retrieval, Joint Embedding, Adversarial Learning, Annotation Efficiency

## 1 INTRODUCTION

Learning robust visual-semantic embeddings is central to the success of many multimedia applications involving multiple modalities such as cross-modal search and data mining [170]. The embedding model aims at encoding and mapping knowledge of multimodal entities into a joint embedding space. The transformation function is typically learned by aligning paired-inputs from two or more distinct domains (*e.g.*, images and natural language descriptions) into the common latent space where the embeddings are close if they are semantically associated or distant if uncorrelated.

Recently, deep neural networks have made significant advancement for learning joint embeddings [49, 93, 105, 140, 210]. Such success is largely attributed to the availability of large-scale human-annotated parallel corpora such as the MS-COCO [116] and Flickr30K [190]

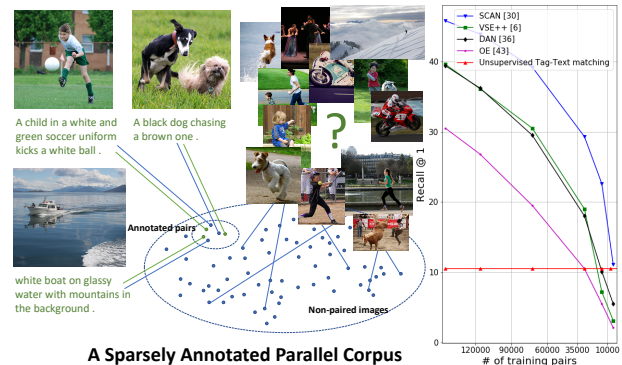


Figure 1: (Left) We consider learning under a sparsely annotated parallel corpus with abundant un-annotated images and limited (image, natural language sentence) pairs. (Right) Performance degradation of state-of-the-art cross-modal retrieval models in the text-to-image retrieval task on Flickr30K. (5 sentences/image.)

datasets. Essentially, there are more than 610,000 and 150,000 annotated image-text pairs in MS-COCO and Flickr30K, respectively. As pointed out in [42], on par to quantity, the annotation diversity is also crucial for downstream tasks. Although models trained with affluent amount of well-annotated image-text pairs can achieve reasonable performance, we observe that the trend does not generalize to more common cases where only a limited amount of parallel annotations are available. As shown in Figure 1, recent VSE models [49, 105, 140, 165] all suffer greater degeneration as annotations become more sparsely available. (See Sec. 4 for experimental details.) Since collecting massive and high-quality human annotations for multimedia corpus is often prohibitively expensive and impractical, there is a strong incentive to designing annotation efficient algorithms to reduce the cost.

In this paper, we deal with the **sparse parallel corpus** scenario (Figure 1) where for cross-modal search and retrieval, a large collection of visual data is available but only a small amount of them are annotated with corresponding text descriptions. We pose an challenging yet rewarding question: *Can we learn satisfactory visual-semantic embedding with a sparse parallel corpus?* Despite some recent progress [97, 107, 159], learning with small amount of parallel data is still challenging and to be developed in urgent need.

A straightforward way to deal with a sparse parallel corpus is to directly utilize the machine generated semantics of the images. In [139], Mithun *et al.* proposed a webly approach to utilize the global tags of the images. However, without handling the inevitable domain gap between the natural language description and

the machine generated tags properly, the visual-semantic embedding learning could be negatively affected, which largely limits the performance.

To circumvent these issues, inspired by the observation in [3] where bottom-up attention over regional objects aligns well with human’s visual system, we propose to utilize “*regional semantics*” which correspond to the regions-of-interest in the un-annotated images and leverage the textual sequences of them to form “pseudo” image-text pairs as the additional weak supervision to conquer the sparsity of image-text annotation. Each regional semantic consists of the category of visual object and its attributes (*e.g. white cat*) which can be automatically extracted with object detection modules [1, 152]. With the inferred regional semantics, we develop a novel method to learn the joint visual-semantic embedding space from both the annotated pairs and the inferred pairs efficiently. To minimize the inherent domain gaps between annotated and un-annotated portion of visual and textual domains, we further impose an attentive alignment with adversarial learning objectives to selectively improve the correlation of semantically close components.

We conduct extensive experiments to quantify the degeneration of current state-of-the-art cross-modal retrieval models in the practical sparse parallel corpus scenario and to show the superiority of the proposed adversarial attentive alignment model for learning visual-semantic embeddings (A3VSE). In terms of reducing annotation effort, in comparison to various recently benchmarks trained with the complete annotations, the proposed model achieves a competitive performance with only 20% of annotations (Recall at 10 > 80.0% for 1K text-to-image and 70.0% for 5K text-to-image retrieval on Flickr30K and MS-COCO, respectively).

In a nutshell, our contributions can be summarized as

- We quantify the impact of learning with common sparse parallel corpora for the state-of-the-art cross-modal retrieval models and shed new insight for annotation efficiency.
- We propose to extract and leverage regional semantics to weakly supervise visual-semantic representation learning.
- We introduce adversarial attentive alignment to deal with multiple heterogeneous domain gaps. The attention mechanism emphasizes the visual or textural informative part to enable effective alignment.
- Experimental results of cross-modal retrieval on the Flickr30k and MS-COCO datasets demonstrate the superiority of our method to the state-of-the-art methods, under the same sparse parallel corpus setting. It is worth noting that, even trained with only 20% of the annotations, our model achieves competitive performance to recent models trained with the complete annotations.

## 2 RELATED WORKS

**Visual-Semantic Embeddings for Cross-Modal Retrieval:** Joint visual-semantic embeddings (VSE) have shown great potential in many multimedia tasks, including cross-modal retrieval [54, 93, 98], visual question answering [5, 64], image captioning [3, 182], multi-modal classification [81], etc. Recently, there are increasing interest in developing system to match natural language descriptions to visual data with VSE [49, 93, 166, 173] for cross-modal retrieval.

In former works, the improvements in VSE are mainly processed on two perspectives: feature learning model and loss function. Various feature learning models have been extensively studied. For the textual feature, the conventional models introduce Fisher vectors [147] for word embeddings [138, 146] as in [54, 99, 171, 173]. Alternatively, recurrent neural networks (RNNs) [72] have been applied in many latest models [49, 85, 86, 88, 93, 94, 105, 140] and Zheng *et al.* suggest a convolutional structure in [210]. For the visual feature, VGG [157] and ResNet [71] models are widely implemented in previous works. Recently, Lee *et al.* [105] proposed to extract regional features from Faster-RCNN model [152]. Attention mechanisms also have been studied in the area [85, 94, 105, 140]. These works learn to select input fragments based on the context from either the same modality [85, 93, 140] or from another modality [105] or both [83]. In [86, 167], additional semantic features has been utilized in a multitask schema. In contrast, in this work, we use image-semantic pair as the weak supervision for learning VSE with sparse corpora.

Most recent works in VSE leverage triplet loss [49, 54, 93, 99, 105, 140, 171, 173]. In [98], Kiros *et al.* proposed to use a triplet ranking loss to penalize the model with individual violations across the negatives. In [171, 173], Wang *et al.* add a within-view neighborhood structure-preserving constraints to further preserve the intra-modal structure. In VSE++ [49], Faghri *et al.* empirically show that emphasizing hard negative examples results in robust joint embeddings. Adversarial objective for cross-modal retrieval is firstly introduced in [167, 184] which narrow down the gap between different modalities by regularization via a domain discriminator. Our work generalize the idea about domain alignment and target on a more common but challenging sparse corpora scenario, where all the above models struggle without a plethora of parallel annotations.

**Learning with Limited Supervision:** Training models with sufficient amount of annotated data could achieve considerable performance for cross-modal retrieval. However, in practice it is difficult to obtain a large amount of well-annotated data [97]. To address this problem, several previous works proposed to utilize web images and their meta data as an auxiliary source of training data [107, 159]. Meanwhile, there are studies focusing on learning with limited supervision. Jiang *et al.* proposed a coupled dictionary learning method to learn the class prototypes that utilize the discriminative information of visual space to improve the less discriminative semantic space in [87]. Tsai *et al.* augmented a typical supervised formulation with unsupervised techniques for learning joint embeddings of visual and textual data in [162]. Although promising performance has been obtained, none of these works consider the sparse parallel corpus setting.

To the best of our knowledge, the most relevant work to ours are [62, 139], where the authors resource meta data and image tags (*i.e.* global semantics) to improve learning of joint embedding space. Our work complements their effort in two perspectives: First, we explore the feasibility of automatic regional semantics as they are more similar to natural language descriptions and leverage them for training improved sequential text encoder. Furthermore, we consider to close the inherent heterogeneous domain gaps with adversarial attentive alignment.

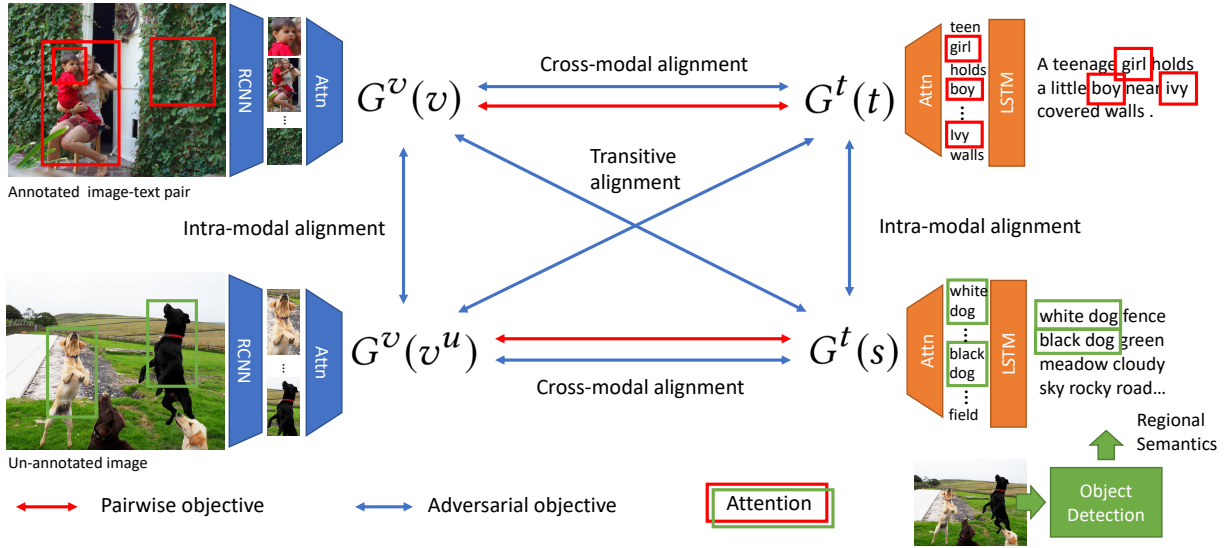


Figure 2: The proposed adversarial attentive alignment model for learning visual-semantic embeddings (A3VSE) for sparsely annotated multimodal corpora. Our model incorporates pseudo “image-text” pairs (illustrated as the bottom image-semantic pair) from the sequence of regional semantics of salient visual objects in un-annotated images. The triplet objectives (colored in red) and adversarial objectives (colored in blue) attend and align semantically correlated instances in the joint embedding space while closing the heterogeneous domain gaps between the annotated/un-annotated portion of visual and textual inputs.

### 3 METHODOLOGY

We consider a common scenario where annotated image-text pairs are sparsely available and un-annotated images are abundant. While manually annotating images with natural language descriptions is expensive, automatically indexing them with semantic tags is relatively efficient [44]. Inspired by the bottom-up approach by [3], instead of resourcing global semantic tags as in [139], we seek to leverage semantics of salient regional objects which aligns well with the natural attention in human’s cognition system to form additional image-semantic pairs for training. However, the inferred regional semantics exhibit clear difference to the natural language descriptions as in the annotated image-text pairs. A judicious way incorporating in these “pseudo” image-semantic pairs across heterogeneous domains for learning visual-semantic embeddings is therefore important.

Figure 2 illustrates the proposed adversarial attentive alignment model for learning visual-semantic embeddings (A3VSE). The proposed model jointly leverages the strong supervision from the annotated image-text pairs and the weak supervision from the inferred image-semantic pairs. Furthermore, A3VSE employs attentive adversarial objectives to selectively align entities from the annotated and un-annotated portion of visual and textual inputs and narrow the domain gaps in between.

#### 3.1 Problem Formulation

Let  $\mathcal{D}^l = \{I_1, \dots, I_{N_l}\}$  be an annotated collection of instances where each instance  $I_i = (v, t)$  consists of the image  $v$  and the corresponding natural language description  $t$ . Let  $\mathcal{D}^u = \{v_1^u, \dots, v_{N_u}^u\}$  denotes the collected but un-annotated images. We name  $\mathcal{D} = \mathcal{D}^l \cup \mathcal{D}^u$  where  $N_l \ll N_u$ , as a **sparse parallel corpus**. We aim

to utilize the un-annotated data  $\mathcal{D}^u$ , together with the annotated data  $\mathcal{D}^l$ , to learn better visual-semantic embeddings.

#### 3.2 Feature Extractors

Let  $F^v$  and  $F^t$  denote the visual feature extractor and the textual feature extractor, respectively. We model  $F^v$  as a fixed object detection model (e.g. Faster RCNN), followed by a trainable fully-connected layer for mapping raw visual features in Faster RCNN into a  $H$ -dimension joint embedding space. On the other hand,  $F^t$  encodes the word tokens in a sentence with a word embedding matrix, followed by a trainable long short-term memory (LSTM) network to model the sequential text inputs. Note that the encoders  $F^v$  and  $F^t$  are shared among  $\mathcal{D}^l$  and  $\mathcal{D}^u$ .

The visual feature of an image  $v$  is encoded as  $V = F^v(v) = [v_1, \dots, v_N] \in \mathbb{R}^{H \times N}$ , where  $N$  is the maximum number of region-of-interest. Similarly, a sentence  $t = [t_1, \dots, t_M]$  is encoded as  $T = F^t(t) = [t_1, \dots, t_M] \in \mathbb{R}^{H \times M}$ , where  $M$  is the maximum sentence length.  $(V_i, T_i)$  represents an annotated feature pair.

For  $v^u \in \mathcal{D}^u$ , we utilize an object detector (Faster RCNN [152]) to extract sequences of regional semantics (as text tokens,  $s = [s_1, \dots, s_M]$ ) and generate image-semantic pairs  $(V_i^u, S_i)$ . The regional semantics are the word tokens of attribute and the class name of the objects detected from an image  $v^u$  (e.g. “blue car”). The detected textual tokens are sorted by their object-wise confidence scores. We concatenate the regional semantics into one sentence, and then encode it as  $S = [s_1, \dots, s_M] \in \mathbb{R}^{H \times M}$  via the shared  $F^t$ .

#### 3.3 Adversarial Attentive Alignment

For learning and aligning instance-wise representation in individual modalities, we apply an attention network which focuses on certain encoded region/ tokens of inputs with respect to the global context

from the same modality. We leverage a  $K$ -head context-aware attention network to capture the interactions between encoded entities and select informative ones for cross-modal alignment.

Given the feature representations (*i.e.* the visual features  $\mathbf{V}$  or the texture features  $\mathbf{T}$ ), the attentive encoder can be written as (we take visual features as an example):

$$E^v(\mathbf{V}) = [\mathbf{W}_0^v \mathbf{V}^\top, \mathbf{W}_1^v \mathbf{V}^\top, \dots, \mathbf{W}_{K-1}^v \mathbf{V}^\top] \quad (1)$$

where

$$W_{ik}^v = \frac{\exp(\lambda_v \alpha_{ik}^v)}{\sum_{i'} \exp(\lambda_v \alpha_{i'k}^v)},$$

$$\alpha_{ik}^v = \tanh(\mathbf{P}_k^v \frac{1}{M^v} \sum_{i'} v_{i'}^\top) \tanh(\mathbf{Q}_k^v v_i).$$

The  $\mathbf{W}_k^v \in \mathbb{R}^{1 \times M^v}$ , and  $\mathbf{P}_k^v, \mathbf{Q}_k^v \in \mathbb{R}^{K' \times H}$ ,  $k \in \{0, 1, \dots, K-1\}$  are the parameters of the attentive encoder  $E^v$ , *i.e.*  $\theta_{v-attn} = \{(\mathbf{W}_k^v, \mathbf{P}_k^v, \mathbf{Q}_k^v) | k \in \{0, 1, \dots, K-1\}\}$ . The  $\lambda_v$  is a constant temperature for the softmax function. The attentive encoder for the textual features (denoted by  $E^t(\mathbf{T})$ ) works the same way but with independent parameters  $\theta_{t-attn} = \{(\mathbf{W}_k^t, \mathbf{P}_k^t, \mathbf{Q}_k^t) | k \in \{0, 1, \dots, K-1\}\}$ . Note that  $E^t$  and  $E^v$  are shared among  $\mathcal{D}_l$  and  $\mathcal{D}_u$ .

Thus, for an image  $v$  or  $v^u$ , the instance-level feature representation can be extracted and selectively encoded through  $G^v = E^v \circ F^v$ . Correspondingly, for the text description  $t$  or  $s$ , the instance-level feature can be achieved by  $G^t = E^t \circ F^t$ . We use  $\theta_v = \{\theta_{v-attn}, \theta_{v-enc}\}$  and  $\theta_t = \{\theta_{t-attn}, \theta_{t-enc}\}$  to denote the trainable parameters of  $G^v$  and  $G^t$ , respectively.

**Triplet Alignment.** For learning the joint embedding, we apply a hinge-based triplet ranking loss with hard negative mining as in [49] to align instance-wise paired visual-textual representations. Let  $(\mathbf{a}, \mathbf{b})$  denotes a sampled image-text or image-semantic pair and  $S(\mathbf{a}, \mathbf{b})$  is the cosine similarity. Let  $\hat{\mathbf{b}} = \operatorname{argmax}_{\mathbf{b}^-} S(\mathbf{a}, \mathbf{b}^-)$  and  $\hat{\mathbf{a}} = \operatorname{argmax}_{\mathbf{a}^-} S(\mathbf{a}^-, \mathbf{b})$  denote the hard negatives in the sampled batch. The triplet objective can be written as:

$$\ell^p(\mathcal{A}, \mathcal{B}; \alpha) = \frac{1}{L} \sum_{i=1}^L \{[\alpha - S(\mathbf{a}_i, \mathbf{b}_i) + S(\mathbf{a}_i, \hat{\mathbf{b}})]_+ + [\alpha - S(\mathbf{a}_i, \mathbf{b}_i) + S(\hat{\mathbf{a}}, \mathbf{b}_i)]_+\}, \quad (2)$$

where  $|\mathcal{A}| = |\mathcal{B}| = L$ ,  $[\cdot]_+ = \max(0, \cdot)$ , and  $\alpha$  is the margin between the similarity of positive pair and that of hard-negative pair. Since annotated image-text pairs sampled from  $\mathcal{D}^l$  are more reliable than image-semantic pairs sampled from  $\mathcal{D}^u$ , we differentiate the strong supervision by the former from the later with a hyper-parameter  $\beta$ . We model the triplet alignment objective as:

$$\ell^{tri} = \beta \ell^p(G^v(v), G^v(t); \alpha_{vt}) + (1 - \beta) \ell^p(G^v(v^u), G^t(s); \alpha_{vs}) \quad (3)$$

A3VSE takes four different types of data, *i.e.*  $\mathbf{V}, \mathbf{T}, \mathbf{V}^u, \mathbf{S}$  which are regarded as samples from four different domains. As shown in Figure 2, we propose using adversarial training to minimize the domain gaps among them. Specifically, we introduce six domain discriminators which are parameterized by  $\theta_{vv^u}, \theta_{ts}, \theta_{vt}, \theta_{v^u s}, \theta_{vs}$ , and  $\theta_{v^u t}$ . On one hand, they are trained to classify samples into correct domains. On the other hand, we employ the gradient reversal layer (GRL) [56] to the reverse the gradients propagated from

these discriminators to update  $G^v$  and  $G^t$  to minimize the domain discrepancy. Such adversarial process can effectively diminish the discrepancy across different domains.

Generally, the adversarial loss for aligning two domains is

$$\ell^d(\mathcal{A}, \mathcal{B}; \theta) = \frac{1}{|\mathcal{A}|} \sum_{i=1}^{|\mathcal{A}|} \log D_\theta(\mathbf{a}_i) + \frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \log(1 - D_\theta(\mathbf{b}_j)) \quad (4)$$

where  $D_\theta$  is the domain discriminator parameterized by  $\theta$ . The  $\mathcal{A} = \{\mathbf{a}_i\}$  and  $\mathcal{B} = \{\mathbf{b}_j\}$  are the mini-batch data sampled from two domains. The instantiations of  $\mathbf{a}$  and  $\mathbf{b}$  can be either two of  $\{G^v(v), G^v(v^u), G^t(t), G^t(s)\}$ . As shown in Figure. 2, we perform three types of alignments, *i.e.* *intra-modal alignment*, *Cross-modal alignment*, and *Transitive alignment*, which are described as follows.

**Intra-modal Alignment** handles the domain gaps between the annotated and un-annotated images, and annotated text descriptions and sequences of regional semantics. Specifically,

$$\ell^{intra} = \lambda_{vv^u} \ell^d(G^v(v), G^v(v^u); \theta_{vv^u}) + \lambda_{ts} \ell^d(G^t(t), G^t(s); \theta_{ts}) \quad (5)$$

**Cross-modal Alignment** aims at aligning the distribution of attended visual and textual features for annotated image-text pairs and inferred image-semantic pairs. That is,

$$\ell^{cross} = \lambda_{vt} \ell^d(G^v(v), G^t(t); \theta_{vt}) + \lambda_{v^u s} \ell^d(G^v(v^u), G^t(s); \theta_{v^u s}) \quad (6)$$

**Transitive Alignment** minimizes the domain gap between annotated images and sequences of regional semantics, and the domain gap between un-annotated images and annotated text descriptions:

$$\ell^{trans} = \lambda_{vs} \ell^d(G^v(v), G^t(s); \theta_{vs}) + \lambda_{v^u t} \ell^d(G^v(v^u), G^t(t); \theta_{v^u t}) \quad (7)$$

The overall adversarial objective for the attentive alignment is:

$$\ell^{adv} = \ell^{intra} + \ell^{cross} + \ell^{trans} \quad (8)$$

And the final objective can be formalized as

$$\ell^{A3VSE} = \ell^{adv} + \ell^{tri} \quad (9)$$

### 3.4 Optimization

**Training and Inference.** A min-max optimization is performed between the domain discriminators and attentive encoders:

$$(\theta_v, \theta_t) = \operatorname{argmin}_{\theta_v, \theta_t} \ell^{A3VSE}(\theta) \quad (10)$$

$$(\theta_{adv}) = \operatorname{argmax}_{\theta_{adv}} \ell^{A3VSE}(\theta),$$

where  $\theta_{adv} \triangleq (\theta_{vt}, \theta_{v^u s}, \theta_{ts}, \theta_{vv^u}, \theta_{vs}, \theta_{v^u t})$ . In each iteration, we sample a mini-batch of  $(v, t)$  from  $\mathcal{D}^l$  and  $(v^u, s)$  from  $\mathcal{D}^u$  then follow the common practice in [56] of adversarial training with GRL to optimize Eq. 9. At the inference stage, we extract the visual embedding for image  $v$  and textual embedding for sentence  $t$  through  $G^v$  and  $G^t$ .

**Discussion.** In A3VSE, attentive encoders and adversarial alignment cooperate to learn satisfactory visual-semantic embeddings. On one hand, attentive encoders emphasize the informative part of the visual regions or textual entities, which helps adversarial

training avoid misalignment and learn more discriminative features; on the other hand, adversarial alignment contributes to the improvement of attention mechanism of the attentive encoders in individual modalities which otherwise may be biased by the less amount of parallel image-text data.

## 4 EXPERIMENT

We perform extensive experiments to confirm the superiority of the proposed A3VSE model over competitive baselines with sparsely annotated multimodal corpora. We evaluate the learned visual-semantic embeddings in cross-modal retrieval tasks on two standard benchmark datasets (Flickr30K [190] and MS-COCO [116]) with the main goal of building an annotation efficient cross-modal retrieval model.

### 4.1 Dataset and Metric

We consider two commonly used benchmark datasets with natural language image descriptions: Flickr30K [190] and MS-COCO [116]. We constrain the amount of image-text annotations available in the training phase as an analogy to real-world scenarios where annotations are typically sparsely available.

**Flickr30K** [190]: There are 31,783 images and 158,915 image-text pairs in the Flickr30K dataset. Five English descriptions are annotated for each image. We start with the standard split defined in [93] with 29,000 training, 1,000 validation, and 1,000 testing images. For learning with limited parallel pairs, we randomly shuffle once and trim the training set into 14,500 (50%), 5,800 (10%), and 2,900 (10%) subset of images. We sample 1, 2, and 5 text descriptions corresponding to those images. The resulting sparse training set is with size 2,900 (2%) to 72,500 (50%) out of 145,000 (100%) training image-text pairs in the original training split. The statistics of the new training splits of sparse Flickr30K can be found in Table 1. The standard validation and the testing are used for model selection and testing. **MS-COCO** [116]: The MS-COCO dataset contains 123,287 images where each image is annotated with five English descriptions. In total, 616,435 image-text pairs are available. We follow the widely used split in [93] to move originally left 30,504 validation images to the training set, resulting a training set of 113,287 training images and 566,435 image-text pairs. We follow the same procedure as performed in Flickr30K and sample 5,664 (5%), 11,382 (10%), and 22,657 (20%) images along with 1, 2, 5 corresponding text descriptions. The statistics and the amount of training pairs can be found in Table 3. We report the testing performance on the whole 5,000 testing set. **Metric**: As in most prior work on cross-modal retrieval tasks [49, 105, 140, 210], we measure rank-based performance by recall at  $K$  ( $R@k$ ). Given a query, recall at  $k$  ( $R@k$ ) calculates the percentage of test instances for which the correct one can be found in the top- $K$  retrieved instances. We report  $R@1$ ,  $R@5$ , and  $R@10$ .

### 4.2 Experimental Setup and Baselines

We focus on the text-to-image retrieval task (searching images with a natural language description as the query) and the image-to-text retrieval task (searching sentences with a query image) with the learned visual-semantic embeddings. We train models under different levels of training sparsity. Model selection and testing are with the full validation and the full testing set, respectively.

For all the baselines, we use their best single model settings and the code from their publicly available Github repositories. Since there are much less paired training instances in sparsely annotated dataset, for fair comparison and in prevention of under-fitting, we either keep the number of (mini-batch) training iterations as 50% iterations of the full dataset or extend the training epoch by 1.2x (for 50% annotations), 2.0x (20% annotations) and 2.5x (10% annotations). Early stopping and learning rate adjustment in the baselines follow the same adjustment if feasible.

**Unsupervised baseline with image-level semantics** We build an unsupervised cross-modal retrieval baseline using *NO* parallel annotations. Image-level semantics (*i.e.*, global semantics) of each image are extracted using pre-trained models from the following datasets: (1) Open Image [103]: 5,000 semantics trained on 9 million images. (2) ImageNet Shuffle [136], 12,073 classes defined in ImageNet. (3) Place365 [211]: 365 visual scene types. (4) Google Sports [95]: 478 sport-related semantics. We remove duplicated semantic concepts, normalize the scores, and then merge them into a 16500-dimension global semantic vector  $s_g$  for each image. Each dimension can be referred to a semantic concept in the original dataset. For example, an “aquarium” in Place365.

For retrieval, we directly match image-level semantics (tags) to text. Specifically, we expand the tokens in a sentence with the synsets defined in WordNet[53] and construct a 16500-dimension  $k$ -hot query vector  $q$ , where  $k$  is the number of matched concepts. The matching score is calculated as  $r = s_g^T q$ .

### 4.3 Implementation Details

We now detail the pre-processing and implementation of the proposed model. To identify and vectorize salient visual objects in images, we use the Faster RCNN model [152] in [3] to detect objects and extract their corresponding visual features  $V \in \mathbb{R}^{36 \times 2048}$ . 36 is the maximum number of ROI in an image and 2,048 is the dimension of the flattened 5-th pooling layer of Faster RCNN [152]. We use raw features without l2 normalization.

For regional semantics in un-annotated images, we use the Faster RCNN model in [3] fine-tuned on Visual Genome [102] to extract English attribute names and class names of the objects detected from an image. Specifically, for every un-annotated image  $v_j^u \in \mathcal{D}_u$ , we generate  $s_j = [s_{j1} || s_{j2} \cdots || s_{j|ROI|}]$  where “||” is concatenation and  $s_k = [\text{Attribute}_k \text{ Class}_k]$  (*e.g.* “blue car”). There are 2,000 detectable objects and attributes. These regional semantics are then sorted by the confidence scores and concatenated as a text sequence. We group the image and the sequence and encode them as an image-semantic pair  $(V^u, S)$ .

In our model, we set the embedding dimension  $H$  to 512. The same dimension is shared by all the context vectors in the attention modules. For text pre-processing, we tokenize, lower-case, truncate maximum sentence length to 57 on MS-COCO and 82 on Flickr30K, and then remove word tokens which appear less than 4 times. Similar to [210], we initialize word embeddings with pre-trained Glove embeddings [146]. All the weights within the network are initialized with Xavier initialization [59]. Other hyper-parameters are set as follows:  $K = 3$ ,  $\alpha_{vt} = 0.2$ ,  $\alpha_{vs} = 0.3$ ,  $\beta = 0.8$ , and  $\gamma = 2/(1 + \exp(-\eta p)) - 1$  as in [56] where  $\eta = 10$  and  $p$  is linearly increased from 0 to 1 in proportional to the training epoch. The hyper-parameters for the adversarial object is set as: Intra-modal

Sparse Flickr30K				Ours (A3VSE)						SCAN [105] (SOTA)					
%	#	%	# Ann	Text-to-Image			Image-to-Text			Text-to-Image			Image-to-Text		
Img	Sent	Ann	Pairs	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
10%	1/5	2%	2,900	20.7	46.0	58.5	27.6	56.2	68.1	2.0	7.2	11.7	5.1	16.0	22.9
10%	2/5	4%	5,800	28.1	55.6	66.9	42.0	69.7	79.0	16.1	35.7	46.5	18.9	39.8	53.9
10%	5/5	10%	14,500	32.0	60.1	71.0	46.8	72.8	80.7	24.6	48.1	59.3	25.9	56.3	70.9
20%	1/5	4%	5,800	29.1	56.4	68.1	43.3	71.0	81.8	17.2	37.5	47.5	21.2	44.4	55.0
20%	2/5	8%	11,600	32.6	61.6	72.3	44.8	72.7	82.8	28.4	54.0	64.6	39.0	68.0	78.6
20%	5/5	20%	29,000	34.9	64.4	73.6	48.4	77.0	85.1	29.3	56.9	68.3	42.1	71.8	81.3
50%	1/5	10%	14,500	36.7	65.1	75.9	51.6	78.7	85.7	29.5	56.3	67.3	40.2	72.2	81.4
50%	2/5	20%	23,200	42.9	70.5	80.3	61.4	83.7	89.4	33.9	61.3	71.4	46.8	75.2	84.5
50%	5/5	50%	72,500	44.5	73.8	83.3	60.9	85.7	91.6	39.2	67.5	77.2	52.6	80.3	87.5

**Table 1: Performance comparison on the 1K testing set of Flickr30K. The models are trained with the sparsely annotated training data as specified in the left column. % *Img* stands for the percentage of training images available compared to original training images in Flickr30K. # *Sent* stands for the number of paired text descriptions available for each image. %/# *Ann* is the percentage/number of annotations used for training compared to the complete training annotations in Flickr30K.**

Model	Text-to-Image			Image-to-Text		
	R@1	R@5	R@10	R@1	R@5	R@10
Flickr30K 0% Ann, 0 pairs						
$s_g$ baseline	10.5	21.5	29.2	12.1	24.0	31.1
Flickr30K 10% Img, 5/5 Sent, 10% Ann, 14,500 pairs						
DPC [210]	8.5	26.0	40.9	11.8	45.5	66.0
DAN [140]	10.1	25.3	42.8	12.2	41.7	64.5
VSE++ [49]	7.2	27.5	40.5	10.5	40.2	62.8
SCAN [105]	24.6	48.1	59.3	25.9	56.3	70.1
<b>Ours (A3VSE)</b>	<b>32.0</b>	<b>60.1</b>	<b>71.0</b>	<b>46.8</b>	<b>73.2</b>	<b>80.7</b>
Flickr30K 50% Img, 2/5 Sent, 20% Ann, 29,000 pairs						
DPC [210]	26.4	53.0	63.9	35.8	68.5	79.7
DAN [140]	26.9	52.3	64.8	37.2	69.9	78.2
VSE++ [49]	27.3	54.5	66.0	33.5	65.2	78.2
SCAN [105]	33.9	61.3	71.4	46.8	75.2	84.5
<b>Ours (A3VSE)</b>	<b>42.9</b>	<b>70.5</b>	<b>80.3</b>	<b>61.4</b>	<b>83.7</b>	<b>89.4</b>
Flickr30K 100% Ann, 145,000 pairs						
DPC [210]	39.1	69.2	80.9	55.6	81.9	89.0
DAN [140]	39.4	69.2	79.1	55.0	81.8	89.5
VSE++ [49]	39.6	70.1	79.8	53.1	82.1	87.5
SCAN [105]	45.8	74.4	83.0	61.8	87.5	93.7
<b>Ours (A3VSE)</b>	<b>49.5</b>	<b>79.5</b>	<b>86.6</b>	<b>65.0</b>	<b>89.2</b>	<b>94.5</b>

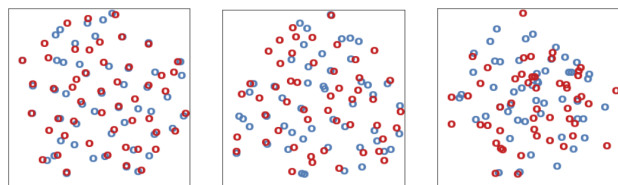
**Table 2: Performance comparison with baselines on two sparse settings in Flickr30K.**

alignments:  $\lambda_{v^u} = 0.2$ ,  $\lambda_{t^s} = 0.1$ ; Cross-modal alignments:  $\lambda_{v^t} = 0.5$ ,  $\lambda_{t^v} = 0.5$ ; Transitive alignments:  $\lambda_{v^u t} = \lambda_{v^s} = 0.3$ .

For training, we train 24 epochs with Adam [96] optimizer. Learning rate is first 0.0005 then 0.00005 after 16th epoch. Models with the greatest summation of recall at 1, 5, 10 in the validation set are selected for testing. Weight decay is set to 0.000001 and gradients larger than 2.0 are clipped. The batch size is 128.

#### 4.4 Results on Sparse Flickr30K

Table 1 shows the testing results with various levels of training sparsity on Flickr30K. Comparing the performance under the same percentage of annotations, the first interesting observation is that generally speaking it is preferred to have diverse images annotated



(a) 20% Img, 5/5 Sent (b) 10% Img, 5/5 Sent (c) No *s* in (b)

**Figure 3: t-SNE visualization of the embedded testing images (blue) and sentences (red) under sparse Flickr30K. Paired ones are expected to be close to each other.**

than annotating a small amount of images with more text descriptions. With the same 10% annotations, it is better to annotate 50% of images with one sentence each than 10% of images with five sentences. These results suggest that regarding data collection and annotation, visual diversity is likely to be more important than textual diversity. Two cases of t-SNE visualization of the learned embedding are shown in Figure 3a and Figure 3b.

Under all sparse training set settings, the proposed model outperforms current state-of-the-art cross-modal retrieval model [105] by a significant margin. Namely, 4.2 to 18.7 in R@1, 6.3 to 38.8 in R@5, and 5.3 to 46.8 in R@10 text-to-image retrieval tasks. Notably, greater improvement over current best model is achieved when less pairwise annotations are available. The improvements converges (but still outperforms) with more annotations available. A similar trend can be observed for the image-to-text retrieval task. These results demonstrate that the proposed A3VSE model can judiciously use regional semantics from un-annotated images for training its encoders and effectively learn the visual-semantic embeddings.

As shown in Table 2, in comparison to other recent models DAN [140], DPC [210], and VSE++ [49], the proposed model significantly outperforms them in all scenarios. In terms of reducing annotation effort, the proposed A3VSE model achieves competitive performance (with the criteria defined as  $R@10 > 80.0\%$ ) trained on only 20% annotations (23,200 pairs).

It is noteworthy that the unsupervised approach with global semantics which use *NO* image-text pairs cannot deliver satisfactory retrieval performance when query with natural language, indicating that there is a clear domain shift between the semantic pool

Sparse MS-COCO				Ours (A3VSE)						SCAN [105] (SOTA)					
%	# Img	# Sent	Ann Pairs	Text-to-Image			Image-to-Text			Text-to-Image			Image-to-Text		
				R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
5%	1/5	1%	5,664	14.2	35.8	48.9	19.2	44.2	57.4	9.0	24.4	35.0	9.5	27.2	38.9
5%	2/5	2%	11,328	16.1	39.5	52.8	22.2	47.8	61.8	12.7	31.6	42.9	11.9	33.1	46.1
5%	5/5	5%	28,320	19.7	44.4	57.7	27.8	55.9	68.8	16.8	40.0	52.6	21.0	47.3	61.2
10%	1/5	2%	11,328	17.7	41.9	54.8	24.6	51.5	63.7	12.7	31.8	43.2	12.8	34.1	48.2
10%	2/5	4%	22,656	20.3	45.5	58.8	26.5	55.6	68.8	17.3	41.5	54.4	22.4	49.7	62.5
10%	5/5	10%	56,640	23.2	50.5	64.1	30.5	60.4	73.1	19.4	44.3	57.3	25.5	53.8	67.6
20%	1/5	4%	22,657	20.0	45.9	59.5	26.9	54.4	67.9	16.3	37.9	50.3	17.8	43.4	57.0
20%	2/5	8%	45,314	24.5	51.8	64.8	32.4	63.0	75.1	20.3	44.5	57.3	24.2	53.7	67.5
20%	5/5	20%	113,287	27.4	56.0	68.9	38.3	68.1	79.3	21.1	45.2	57.8	24.2	54.8	68.6

Table 3: Performance comparison on the 5K testing set of MS-COCO.

Model	Text-to-Image			Image-to-Text		
	R@1	R@5	R@10	R@1	R@5	R@10
MS-COCO 0% Ann, 0 pairs						
$s_g$ baseline	7.5	16.8	23.2	8.8	15.0	22.8
MS-COCO 10% Img, 1/5 Sent, 2% Ann, 11,328 pairs						
DPC [210]	8.1	28.3	38.0	10.5	30.8	41.0
DAN [140]	8.8	28.3	37.1	11.1	30.1	42.5
VSE++ [49]	8.5	27.6	36.5	10.7	30.2	44.5
SCAN [105]	12.7	31.8	43.2	12.8	34.1	48.2
<b>Ours (A3VSE)</b>	<b>17.7</b>	<b>41.9</b>	<b>54.8</b>	<b>24.6</b>	<b>51.5</b>	<b>63.7</b>
MS-COCO 50% Img, 2/5 Sent, 20% Ann, 113,287 pairs						
DPC [210]	19.1	41.0	55.5	20.5	45.1	60.2
DAN [140]	19.5	40.8	54.0	20.7	47.7	61.7
VSE++ [49]	19.5	41.2	56.5	21.5	48.5	63.5
SCAN [105]	22.3	47.5	60.2	25.5	56.1	70.5
<b>Ours (A3VSE)</b>	<b>28.2</b>	<b>57.9</b>	<b>70.6</b>	<b>38.4</b>	<b>69.5</b>	<b>81.1</b>
MS-COCO 100% Ann, 566,435 pairs						
DPC [210]	25.3	53.4	66.4	41.2	70.5	81.1
DAN [140]	29.8	58.8	70.0	40.8	70.0	79.8
VSE++ [49]	30.3	56.0	72.4	41.3	69.5	81.2
SCAN [105]	34.4	63.7	75.7	46.4	77.4	87.2
<b>Ours (A3VSE)</b>	<b>39.0</b>	<b>68.0</b>	<b>80.1</b>	<b>49.3</b>	<b>81.1</b>	<b>90.2</b>

Table 4: Performance comparison with baselines on two sparse settings in MS-COCO.

of current image classification/ tagging models and the natural language queries. A similar phenomena is observed in our ablation study. Moreover, from the crossover of 10.5 R@1 in Figure 1 (right), the unsupervised global semantics from external classification datasets is worth as many as 14,000 image-text annotation pairs for the recent cross-modal retrieval models. Notably, A3VSE achieves 29.1 R@1 even trained with only 5,800 pairs.

#### 4.5 Results on Sparse MS-COCO

Table 3 shows the results on the harder 5K testing set of MS-COCO. We sample 5%, 10%, 20% of images in MS-COCO to keep the number training of pairs more comparable to Flickr30K. The proposed model delivers the best performance on most metrics under all sparsity settings. For text-to-image retrieval, it outperforms SCAN [105] by 2.9 to 6.3 in R@1, 4.0 to 11.4 in R@5, and 4.4 to 13.9 in R@10. Similar trend can be observed in image-to-text retrieval task. The

Flickr30K 10% Img 5/5 Sent, 10% Ann, 14,500 pairs						
Model	Text-to-Image			Image-to-Text		
	R@1	R@5	R@10	R@1	R@5	R@10
No $s$	23.4	47.9	58.2	26.5	58.1	71.5
Swap $s$ with $s_g$	29.0	56.3	67.2	40.5	67.4	77.6
$s$ , without attention	23.8	50.1	62.7	35.8	64.3	75.1
$s$ , without $L_{adv}$	30.9	58.5	69.0	43.8	70.9	79.5
Without $\ell^{intra}$	31.8	59.6	<b>71.0</b>	44.8	72.5	<b>80.8</b>
Without $\ell^{cross}$	31.3	59.2	70.5	45.2	71.8	80.1
Without $\ell^{trans}$	31.5	59.7	70.9	46.1	71.8	80.3
Full model	<b>32.0</b>	<b>60.1</b>	<b>71.0</b>	<b>46.8</b>	<b>72.8</b>	80.7

Table 5: Ablation study of the proposed model

comparison with other recently published models is shown in Table 4 where the proposed model achieves the best performance in all sparse corpus scenarios.

Despite using only 20% of image-text annotations, the proposed model still achieves competitive performance (with the criteria defined as  $R@10 > 70.0\%$ ) in the more challenging 5K testing set in MS-COCO. More than 80% of annotation effort for the image-text pairs could potentially be relieved. Based on the quantitative results on multiple datasets, we validate the superiority and the annotation efficiency of the proposed A3VSE model.

#### 4.6 Ablation Study

To quantify the contribution from individual components, we conduct ablation studies evaluating the cross-modal retrieval performance with models trained with 10% of images and 5/5 corresponding text descriptions (10% annotations) in Flickr30K. In each experiment, we remove one or change component of concern to quantify its relative importance. The larger the drop implies that the component is more important. For the experiment without semantics ( $s$ ), we remove all the regional semantics from the input and show the performance of the vanilla model. Then we swap the sequence of regional semantics with global semantics  $s_g$  and encode global semantics (can be viewed as image-level tags after applying a 0.3 threshold) with the shared word embedding matrix. For the internal modules and adversarial objectives, we either remove the attention layer with mean pooling over encoded visual/textual entities as the final instance-level representation, or we purge an adversarial objective from Eq. 9 during the training phase.

Table 5 shows the results of the ablation study. We observe that while global semantics boost model performance from the vanilla





(a) 50% Img, 2/5 Sent, 23200 pairs

(b) 10% Img, 2/5 Sent, 5,800 pairs

(c) Failures of (b)

**Figure 4: Qualitative examples of the proposed A3VSE model in text-to-image retrieval task (the upper two rows) and image-to-text retrieval task (the bottom row) on Flickr30K.**

model, the regional semantics is the better choice even if they have a relatively small vocabulary size (1,104 versus 1,576) for the un-annotated images in sparse Flickr30K. The visualization of learned embeddings in Figure 3b and Figure 3c double confirms the difference. One possible explanation for this phenomena is that regional semantics are more similar to natural language descriptions. We observe that the distribution of vocabulary is closer (13.1% Intersection over Union (IoU)) between the natural language queries and the regional semantics than the global semantics (9.8% IoU). For instance, in a natural language description, people tend to describe an image with “frog” or “dog” rather than the detected global semantics “Amphibian” and “havanese”.

Additionally, the attentive adversarial learning with domain discriminators plays an important role for closing the domain gaps between annotated and un-annotated inputs, delivers improved performance over models without adversarial objectives. However, we observe small variants among the best metrics over various configurations, suggesting that a careful hyper-parameter tuning may be required to achieve the optimal performance. We leave the robust automatic tuning for aligning multiple heterogeneous domains as our future work.

## 4.7 Qualitative Results

Figure 4 illustrates sampled qualitative testing results in the image-to-text and text-to-image retrieval tasks on sparse Flickr30K. The top two rows show the top four retrieved images given the natural language query above. The one and only one correct image is marked in green or red if rank > 10. The image-to-text retrieval results are depicted in the bottom row. We list the top five retrieved sentences and the corresponding query image. The correct sentences (up to five) are colored in green otherwise red.

In most cases the proposed model generates satisfactory results. As less parallel image-text pairs are available for training, we observe performance degeneration. For the failure cases, as expected, we observe that many failures result from out-of-vocabulary words (e.g. “amplifier” and “harp”) in the sentences.

## 5 CONCLUSION

To reduce expensive human annotation cost, we have presented a novel annotation efficient A3VSE model for learning improved

visual-semantic embeddings (VSEs) with sparsely annotated multimodal corpora. The proposed model jointly leverages strong supervision from image-text pairs and weak supervision from image-semantic pairs where the regional semantics are extracted from the un-annotated image collection. To further unify the heterogeneous inputs in the joint embedding space, our model employs attention-enhanced adversarial objectives to model intra-modal, cross-modal, and transitive alignment to selectively align annotated and un-annotated portion of visual and textual inputs.

In sparse Flickr30K and MS-COCO, the proposed model consistently and significantly outperforms recent competitive baselines. In comparison to global semantic tags, we have shown that regional semantics are more feasible for learning VSEs under sparsity. With regard to reducing annotation effort, we have presents insights towards efficient annotation collection and utilization. We have demonstrated that nearly 80% of the annotations can be reduced with the proposed model while achieving competitive results to recent models trained with the complete annotations.

## ACKNOWLEDGEMENT

This research is supported by DARPA grant FA8750-18-2-0018 and FA8750-19-2-0501 funded under the AIDA program and the LwLL program. It is also supported by the Intelligence Advanced Research Projects Activity (IARPA) via Department of Interior/Interior Business Center (DOI/IBC) contract number D17PC00340.

## REFERENCES

- [1] Waleed Abdulla. 2017. Mask R-CNN for object detection and instance segmentation on Keras and TensorFlow. [https://github.com/matterport/Mask\\_RCNN](https://github.com/matterport/Mask_RCNN).
- [2] Dharma P. Agrawal, Brij Bhooshan Gupta, Haoxiang Wang, Xiaojun Chang, Shingo Yamaguchi, and Gregorio Martínez Pérez. 2018. Guest Editorial Deep Learning Models for Industry Informatics. *IEEE Trans. Ind. Informatics* 14, 7 (2018), 3166–3169.
- [3] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering. In *CVPR*.
- [4] Galen Andrew, Raman Arora, Jeff Bilmes, and Karen Livescu. 2013. Deep canonical correlation analysis. In *International Conference on Machine Learning*. 1247–1255.
- [5] Stanislaw Antol, Aishwarya Agrawal, Jiaseen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. 2015. VQA: Visual Question Answering. In *International Conference on Computer Vision (ICCV)*.
- [6] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014).

- [7] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural Machine Translation by Jointly Learning to Align and Translate. *CoRR* abs/1409.0473 (2014).
- [8] Xiangpin Bai, Lei Zhu, Cheng Liang, Jingjing Li, Xiushan Nie, and Xiaojun Chang. 2020. Multi-view feature selection via Nonnegative Structured Graph Learning. *Neurocomputing* 387 (2020), 110–122.
- [9] Ozan Caglayan, Walid Aransa, Adrien Bardet, Mercedes Garcia-Martinez, Fethi Bougares, Loïc Barrault, Marc Masana, Luis Herranz, and Joost Van de Weijer. 2017. LIUM-CVC Submissions for WMT17 Multimodal Translation Task. In *SECOND CONFERENCE ON MACHINE TRANSLATION*, Vol. 2. 432–439.
- [10] Iacer Calixto and Qun Liu. 2017. Incorporating Global Visual Features into Attention-based Neural Machine Translation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. 992–1003.
- [11] Xiaojun Chang, Po-Yao Huang, Yi-Dong Shen, Xiaodan Liang, Yi Yang, and Alexander G. Hauptmann. 2018. RCAA: Relational Context-Aware Agents for Person Search. In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part IX*. 86–102.
- [12] Xiaojun Chang, Xiaodan Liang, Yan Yan, and Liqiang Nie. 2020. Guest editorial: Image/video understanding and analysis. *Pattern Recognit. Lett.* 130 (2020), 1–3.
- [13] Xiaojun Chang, Wenhe Liu, Po-Yao Huang, Changlin Li, Fengda Zhu, Mingfei Han, Mingjie Li, Mengyuan Ma, Siyi Hu, Guoliang Kang, Junwei Liang, Liangkue Gui, Lijun Yu, Yijun Qian, Jing Wen, and Alexander G. Hauptmann. 2019. MMVG-INF-Etrol@TRECVID 2019: Activities in Extended Video. In *2019 TREC Video Retrieval Evaluation, TRECVID 2019, Gaithersburg, MD, USA, November 12-13, 2019*.
- [14] Xiaojun Chang, Zhigang Ma, Ming Lin, Yi Yang, and Alexander G. Hauptmann. 2017. Feature Interaction Augmented Sparse Learning for Fast Kinect Motion Detection. *IEEE Trans. Image Process.* 26, 8 (2017), 3911–3920.
- [15] Xiaojun Chang, Zhigang Ma, Yi Yang, Zhiqiang Zeng, and Alexander G. Hauptmann. 2017. Bi-Level Semantic Representation Analysis for Multimedia Event Detection. *IEEE Trans. Cybern.* 47, 5 (2017), 1180–1197.
- [16] Xiaojun Chang, Feiping Nie, Zhigang Ma, Yi Yang, and Xiaofang Zhou. 2015. A Convex Formulation for Spectral Shrunk Clustering. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, January 25-30, 2015, Austin, Texas, USA*. 2532–2538.
- [17] Xiaojun Chang, Feiping Nie, Sen Wang, Yi Yang, Xiaofang Zhou, and Chengqi Zhang. 2016. Compound Rank-k Projections for Bilinear Analysis. *IEEE Trans. Neural Networks Learn. Syst.* 27, 7 (2016), 1502–1513.
- [18] Xiaojun Chang, Feiping Nie, Yi Yang, and Heng Huang. 2014. A Convex Formulation for Semi-Supervised Multi-Label Feature Selection. In *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence, July 27 -31, 2014, Québec City, Québec, Canada*. 1171–1177.
- [19] Xiaojun Chang, Feiping Nie, Yi Yang, Chengqi Zhang, and Heng Huang. 2016. Convex Sparse PCA for Unsupervised Feature Learning. *ACM Trans. Knowl. Discov. Data* 11, 1 (2016), 3:1–3:16.
- [20] Xiaojun Chang, Haoquan Shen, Sen Wang, Jiajun Liu, and Xue Li. 2014. Semi-supervised Feature Analysis for Multimedia Annotation by Mining Label Correlation. In *Advances in Knowledge Discovery and Data Mining - 18th Pacific-Asia Conference, PAKDD 2014, Tainan, Taiwan, May 13-16, 2014. Proceedings, Part II*. 74–85.
- [21] Xiaojun Chang, Yan Yan, and Liqiang Nie. 2018. Guest Editorial: Semantic Concept Discovery in MM Data. *Multim. Tools Appl.* 77, 3 (2018), 2945–2946.
- [22] Xiaojun Chang and Yi Yang. 2017. Semisupervised Feature Analysis by Mining Correlations Among Multiple Tasks. *IEEE Trans. Neural Networks Learn. Syst.* 28, 10 (2017), 2294–2305.
- [23] Xiaojun Chang, Yi Yang, Alexander G. Hauptmann, Eric P. Xing, and Yaoliang Yu. 2015. Semantic Concept Discovery for Large-Scale Zero-Shot Event Detection. In *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina, July 25-31, 2015*. 2234–2240.
- [24] Xiaojun Chang, Yi Yang, Guodong Long, Chengqi Zhang, and Alexander G. Hauptmann. 2016. Dynamic Concept Composition for Zero-Example Event Detection. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA*. 3464–3470.
- [25] Xiaojun Chang, Yi Yang, Eric P. Xing, and Yaoliang Yu. 2015. Complex Event Detection using Semantic Saliency and Nearly-Isotonic SVM. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*. 1348–1357.
- [26] Xiaojun Chang, Yaoliang Yu, and Yi Yang. 2017. Robust Top-k Multiclass SVM for Visual Category Recognition. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Halifax, NS, Canada, August 13 - 17, 2017*. 75–83.
- [27] Xiaojun Chang, Yaoliang Yu, Yi Yang, and Alexander G. Hauptmann. 2015. Searching Persuasively: Joint Event Detection and Evidence Recounting with Limited Supervision. In *Proceedings of the 23rd Annual ACM Conference on Multimedia Conference, MM '15, Brisbane, Australia, October 26 - 30, 2015*. 581–590.
- [28] Xiaojun Chang, Yaoliang Yu, Yi Yang, and Eric P. Xing. 2016. They are Not Equally Reliable: Semantic Event Search Using Differentiated Concept Classifiers. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*. 1884–1893.
- [29] Xiaojun Chang, Yaoliang Yu, Yi Yang, and Eric P. Xing. 2017. Semantic Pooling for Complex Event Analysis in Untrimmed Videos. *IEEE Trans. Pattern Anal. Mach. Intell.* 39, 8 (2017), 1617–1632.
- [30] Yan-shuo Chang, Feiping Nie, Zhihui Li, Xiaojun Chang, and Heng Huang. 2017. Refined Spectral Clustering via Embedded Label Propagation. *Neural Comput.* 29, 12 (2017).
- [31] Kaixuan Chen, Lina Yao, Dalin Zhang, Xiaojun Chang, Guodong Long, and Sen Wang. 2019. Distributionally Robust Semi-Supervised Learning for People-Centric Sensing. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*. 3321–3328.
- [32] Kaixuan Chen, Lina Yao, Dalin Zhang, Xianzhi Wang, Xiaojun Chang, and Feiping Nie. 2020. A Semisupervised Recurrent Convolutional Attention Model for Human Activity Recognition. *IEEE Trans. Neural Networks Learn. Syst.* 31, 5 (2020), 1747–1756.
- [33] Xiaojun Chen, Guowen Yuan, Wenting Wang, Feiping Nie, Xiaojun Chang, and Joshua Zhexue Huang. 2018. Local Adaptive Projection Framework for Feature Selection of Labeled and Unlabeled Data. *IEEE Trans. Neural Networks Learn. Syst.* 29, 12 (2018), 6362–6373.
- [34] De Cheng, Xiaojun Chang, Li Liu, Alexander G. Hauptmann, Yihong Gong, and Nanning Zheng. 2017. Discriminative Dictionary Learning With Ranking Metric Embedded for Person Re-Identification. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*. 964–970.
- [35] De Cheng, Yihong Gong, Xiaojun Chang, Weiwei Shi, Alexander G. Hauptmann, and Nanning Zheng. 2018. Deep feature learning via structured graph Laplacian embedding for person re-identification. *Pattern Recognit.* 82 (2018), 94–104.
- [36] Xuelian Cheng, Yiran Zhong, Mehrtash Harandi, Yuchao Dai, Xiaojun Chang, Tom Drummond, Hongdong Li, and Zongyuan Ge. 2020. Hierarchical Neural Architecture Search for Deep Stereo Matching. *CoRR* abs/2010.13501 (2020).
- [37] Xuelian Cheng, Yiran Zhong, Mehrtash Harandi, Yuchao Dai, Xiaojun Chang, Hongdong Li, Tom Drummond, and Zongyuan Ge. 2020. Hierarchical Neural Architecture Search for Deep Stereo Matching. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- [38] Zhiyong Cheng, Xiaojun Chang, Lei Zhu, Rose Catherine Kanjirathinkal, and Mohan S. Kankanahalli. 2019. MMALFM: Explainable Recommendation by Leveraging Reviews and Images. *ACM Trans. Inf. Syst.* 37, 2 (2019), 16:1–16:28.
- [39] Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the Properties of Neural Machine Translation: Encoder–Decoder Approaches. *Syntax, Semantics and Structure in Statistical Translation* (2014), 103.
- [40] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 1724–1734.
- [41] Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*. 1724–1734.
- [42] Bo Dai, Sanja Fidler, Raquel Urtasun, and Dahua Lin. 2017. Towards Diverse and Natural Image Descriptions via a Conditional GAN. In *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2989–2998.
- [43] Bo Dai and Dahua Lin. 2017. Contrastive learning for image captioning. In *Advances in Neural Information Processing Systems*. 898–907.
- [44] Jia Deng, Alexander C Berg, and Li Fei-Fei. 2011. Hierarchical semantic indexing for large scale image retrieval. In *CVPR 2011*. IEEE, 785–792.
- [45] Michael Denkowski and Alon Lavie. 2014. Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the ninth workshop on statistical machine translation*. 376–380.
- [46] Jianfeng Dong. 2017. Cross-media Relevance Computation for Multimedia Retrieval. In *Proceedings of the 25th ACM international conference on Multimedia*. ACM, 831–835.
- [47] Aviv Eisenschat and Lior Wolf. 2017. Linking Image and Text with 2-Way Nets. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*. IEEE, 1855–1865.
- [48] Desmond Elliott, Stella Frank, Khalil Sima'an, and Lucia Specia. 2016. Multi30K: Multilingual English-German Image Descriptions. In *Proceedings of the 5th Workshop on Vision and Language*. Association for Computational Linguistics, 70–74. <https://doi.org/10.18653/v1/W16-3210>

- [49] Fartash Faghri, David J Fleet, Jamie Ryan Kiros, and Sanja Fidler. 2018. VSE++: Improving Visual-Semantic Embeddings with Hard Negatives. (2018). <https://github.com/fartashf/vsepp>
- [50] Hehe Fan, Xiaojun Chang, De Cheng, Yi Yang, Dong Xu, and Alexander G. Hauptmann. 2017. Complex Event Detection by Identifying Reliable Shots from Untrimmed Videos. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*. 736–744.
- [51] Mingyu Fan, Xiaojun Chang, and Dacheng Tao. 2017. Structure Regularized Unsupervised Discriminant Feature Analysis. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*. 1870–1876.
- [52] Mingyu Fan, Xiaojun Chang, Xiaoqin Zhang, Di Wang, and Liang Du. 2017. Top-k Supervise Feature Selection via ADMM for Integer Programming. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*. 1646–1653.
- [53] Christiane Fellbaum. 1998. *WordNet: An Electronic Lexical Database*. Bradford Books.
- [54] Andrea Frome, Gregory S. Corrado, Jonathon Shlens, Samy Bengio, Jeffrey Dean, Marc’Aurelio Ranzato, and Tomas Mikolov. 2013. DeViSE: A Deep Visual-Semantic Embedding Model. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States*. 2121–2129.
- [55] Yaroslav Ganin and Victor Lempitsky. 2014. Unsupervised domain adaptation by backpropagation. *arXiv preprint arXiv:1409.7495* (2014).
- [56] Yaroslav Ganin and Victor Lempitsky. 2015. Unsupervised Domain Adaptation by Backpropagation. In *International Conference on Machine Learning*. 1180–1189.
- [57] Zongyuan Ge, Dwarikanath Mahapatra, Xiaojun Chang, Zetao Chen, Lianhua Chi, and Huimin Lu. 2020. Improving multi-label chest X-ray disease diagnosis by exploiting disease and health labels dependencies. *Multim. Tools Appl.* 79, 21-22 (2020), 14889–14902.
- [58] Spandana Gella, Rico Senrich, Frank Keller, and Mirella Lapata. 2017. Image Pivoting for Learning Multilingual Multimodal Representations. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2839–2845. <https://doi.org/10.18653/v1/D17-1303>
- [59] Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 249–256.
- [60] Chen Gong, Xiaojun Chang, Meng Fang, and Jian Yang. 2018. Teaching Semi-Supervised Classifier via Generalized Distillation. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*. 2156–2162.
- [61] Chen Gong, Dacheng Tao, Xiaojun Chang, and Jian Yang. 2019. Ensemble Teaching for Hybrid Label Propagation. *IEEE Trans. Cybern.* 49, 2 (2019), 388–402.
- [62] Yunchao Gong, Liwei Wang, Micah Hodosh, Julia Hockenmaier, and Svetlana Lazebnik. 2014. Improving image-sentence embeddings using large weakly annotated photo collections. In *European conference on computer vision*. Springer, 529–545.
- [63] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *Advances in neural information processing systems*. 2672–2680.
- [64] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [65] Jiuxiang Gu, Jianfei Cai, Shafiq Joty, Li Niu, and Gang Wang. [n. d.]. Look, Imagine and Match: Improving Textual-Visual Cross-Modal Retrieval with Generative Models.
- [66] Liangke Gui, Xiaodan Liang, Xiaojun Chang, and Alexander G. Hauptmann. 2018. Adaptive Context-aware Reinforced Agent for Handwritten Text Recognition. In *British Machine Vision Conference 2018, BMVC 2018, Newcastle, UK, September 3-6, 2018*. 207.
- [67] Junwei Han, Le Yang, Dingwen Zhang, Xiaojun Chang, and Xiaodan Liang. 2018. Reinforcement Cutting-Agent Learning for Video Object Segmentation. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*. 9080–9089.
- [68] Longfei Han, Dingwen Zhang, Dong Huang, Xiaojun Chang, Jun Ren, Senlin Luo, and Junwei Han. 2017. Self-paced Mixture of Regressions. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*. 1816–1822.
- [69] Mingfei Han, Yali Wang, Xiaojun Chang, and Yu Qiao. 2020. Mining Inter-Video Proposal Relations for Video Object Detection. In *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XXI*. 431–446.
- [70] Di He, Yingce Xia, Tao Qin, Liwei Wang, Nenghai Yu, Tieyan Liu, and Wei-Ying Ma. 2016. Dual learning for machine translation. In *Advances in Neural Information Processing Systems*. 820–828.
- [71] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [72] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [73] Micah Hodosh, Peter Young, and Julia Hockenmaier. 2013. Framing image description as a ranking task: Data, models and evaluation metrics. *Journal of Artificial Intelligence Research* 47 (2013), 853–899.
- [74] Siyi Hu and Xiaojun Chang. 2020. Multi-view Drone-based Geo-localization via Style and Spatial Alignment. *CoRR abs/2006.13681* (2020).
- [75] Siyi Hu, Fengda Zhu, Xiaojun Chang, and Xiaodan Liang. 2021. UPDeT: Universal Multi-agent Reinforcement Learning via Policy Decoupling with Transformers. *CoRR abs/2101.08001* (2021).
- [76] Po-Yao Huang, Xiaojun Chang, and Alexander G. Hauptmann. 2019. Multi-Head Attention with Diversity for Learning Grounded Multilingual Multimodal Representations. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*. 1461–1467.
- [77] Po-Yao Huang, Xiaojun Chang, Alexander G. Hauptmann, and Eduard H. Hovy. 2020. Forward and Backward Multimodal NMT for Improved Monolingual and Multilingual Cross-Modal Retrieval. In *Proceedings of the 2020 International Conference on Multimedia Retrieval, ICMR 2020, Dublin, Ireland, June 8-11, 2020*. 53–62.
- [78] Po-Yao Huang, Junjie Hu, Xiaojun Chang, and Alexander G. Hauptmann. 2020. Unsupervised Multimodal Neural Machine Translation with Pseudo Visual Pivoting. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*. 8226–8237.
- [79] Po-Yao Huang, Guoliang Kang, Wenhe Liu, Xiaojun Chang, and Alexander G. Hauptmann. 2019. Annotation Efficient Cross-Modal Retrieval with Adversarial Attentive Alignment. In *Proceedings of the 27th ACM International Conference on Multimedia, MM 2019, Nice, France, October 21-25, 2019*. 1758–1767.
- [80] Po-Yao Huang, Vaibhav, Xiaojun Chang, and Alexander G. Hauptmann. 2019. Improving What Cross-Modal Retrieval Models Learn through Object-Oriented Inter- and Intra-Modal Attention Networks. In *Proceedings of the 2019 International Conference on Multimedia Retrieval, ICMR 2019, Ottawa, ON, Canada, June 10-13, 2019*. 244–252.
- [81] Po-Yao Huang, Junwei Liang, Jean-Baptiste Lamare, and Alexander G. Hauptmann. 2018. Multimodal Filtering of Social Media for Temporal Monitoring and Event Analysis. In *Proceedings of the 2018 ACM on International Conference on Multimedia Retrieval (ICMR ’18)*. ACM, New York, NY, USA, 450–457. <https://doi.org/10.1145/3206025.3206079>
- [82] Po-Yao Huang, Frederick Liu, Sz-Rung Shiang, Jean Oh, and Chris Dyer. 2016. Attention-based multimodal neural machine translation. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, Vol. 2. 639–645.
- [83] Po-Yao Huang, Vaibhav, Xiaojun Chang, and Alexander G. Hauptmann. 2019. Improving What Cross-Modal Retrieval Models Learn Through Object-Oriented Inter- and Intra-Modal Attention Networks. In *Proceedings of the 2019 International Conference on Multimedia Retrieval (ICMR ’19)*. ACM, New York, NY, USA, 244–252. <https://doi.org/10.1145/3323873.3325043>
- [84] Po-Yao Huang, Ye Yuan, Zhenzhong Lan, Lu Jiang, and Alexander G Hauptmann. 2017. Video Representation Learning and Latent Concept Mining for Large-scale Multi-label Video Classification. *arXiv preprint arXiv:1707.01408* (2017).
- [85] Yan Huang, Wei Wang, and Liang Wang. 2017. Instance-Aware Image and Sentence Matching with Selective Multimodal LSTM. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 7254–7262.
- [86] Yan Huang, Qi Wu, and Liang Wang. 2017. Learning semantic concepts and order for image and sentence matching. *arXiv preprint arXiv:1712.02036* (2017).
- [87] Huajie Jiang, Ruiping Wang, Shiguang Shan, and Xilin Chen. 2018. Learning class prototypes via structure alignment for zero-shot recognition. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 118–134.
- [88] Xinyang Jiang, Fei Wu, Xi Li, Zhou Zhao, Weiming Lu, Siliang Tang, and Yueting Zhuang. 2015. Deep compositional cross-modal learning to rank via local-global alignment. In *Proceedings of the 23rd ACM international conference on Multimedia*. ACM, 69–78.
- [89] Lukasz Kaiser and Samy Bengio. 2016. Can Active Memory Replace Attention?. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*. 3774–3782.
- [90] Nal Kalchbrenner and Phil Blunsom. 2013. Recurrent Continuous Translation Models. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, EMNLP 2013, 18-21 October 2013, Grand Hyatt Seattle, Seattle, Washington, USA, A meeting of SIGDAT, a Special Interest Group of the ACL*. 1700–1709.

- [91] Nal Kalchbrenner, Lasse Espeholt, Karen Simonyan, Aaron van den Oord, Alex Graves, and Koray Kavukcuoglu. 2016. Neural machine translation in linear time. *arXiv preprint arXiv:1610.10099* (2016).
- [92] Nal Kalchbrenner, Lasse Espeholt, Karen Simonyan, Aaron van den Oord, Alex Graves, and Koray Kavukcuoglu. 2016. Neural Machine Translation in Linear Time. *CoRR abs/1610.10099* (2016).
- [93] Andrej Karpathy and Li Fei-Fei. 2015. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3128–3137.
- [94] Andrej Karpathy, Armand Joulin, and Li Fei-Fei. 2014. Deep fragment embeddings for bidirectional image sentence mapping. In *Advances in neural information processing systems*. 1889–1897.
- [95] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. 2014. Large-scale Video Classification with Convolutional Neural Networks. In *CVPR*.
- [96] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [97] Ryan Kiros, Ruslan Salakhutdinov, and Richard S Zemel. 2014. Unifying visual-semantic embeddings with multimodal neural language models. *arXiv preprint arXiv:1411.2539* (2014).
- [98] Ryan Kiros, Ruslan Salakhutdinov, and Richard S. Zemel. 2014. Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models. *NIPS Workshop* (2014).
- [99] Benjamin Klein, Guy Lev, Gil Sadeh, and Lior Wolf. 2015. Associating neural word embeddings with deep image representations using Fisher Vectors. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*. 4437–4446.
- [100] Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander M. Rush. 2017. OpenNMT: Open-Source Toolkit for Neural Machine Translation. In *Proc. ACL*. <https://doi.org/10.18653/v1/P17-4012>
- [101] Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, et al. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th annual meeting of the ACL on interactive poster and demonstration sessions*. Association for Computational Linguistics, 177–180.
- [102] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannic Kalantidis, Li-Jia Li, David A Shamma, Michael Bernstein, and Li Fei-Fei. 2016. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations. <https://arxiv.org/abs/1602.07332>
- [103] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Tom Duerig, et al. 2018. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *arXiv preprint arXiv:1811.00982* (2018).
- [104] Patrik Lambert, Holger Schwenk, Christophe Servan, and Sadaf Abdul-Rauf. 2011. Investigations on translation model adaptation using monolingual data. In *Proceedings of the Sixth Workshop on Statistical Machine Translation*. Association for Computational Linguistics, 284–293.
- [105] Kuang-Huei Lee, Xi Chen, Gang Hua, Houdong Hu, and Xiaodong He. 2018. Stacked Cross Attention for Image-Text Matching. *arXiv preprint arXiv:1803.08024* (2018).
- [106] Changlin Li, Jiefeng Peng, Liuchun Yuan, Guangrun Wang, Xiaodan Liang, Liang Lin, and Xiaoju Chang. 2020. Block-Wisely Supervised Neural Architecture Search With Knowledge Distillation. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*. 1986–1995.
- [107] Junnan Li, Yongkang Wong, Qi Zhao, and Mohan S Kankanhalli. 2017. Attention transfer from web images for video recognition. In *Proceedings of the 25th ACM international conference on Multimedia*. ACM, 1–9.
- [108] Mingjie Li, Fuyu Wang, Xiaoju Chang, and Xiaodan Liang. 2020. Auxiliary Signal-Guided Knowledge Encoder-Decoder for Medical Report Generation. *CoRR abs/2006.03744* (2020).
- [109] Zhihui Li, Xiaoju Chang, Lina Yao, Shirui Pan, Zongyuan Ge, and Huaxiang Zhang. 2020. Grounding Visual Concepts for Zero-Shot Event Detection and Event Captioning. In *KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020*. 297–305.
- [110] Zhihui Li, Wenhe Liu, Xiaoju Chang, Lina Yao, Mahesh Prakash, and Huaxiang Zhang. 2019. Domain-Aware Unsupervised Cross-dataset Person Re-identification. In *Advanced Data Mining and Applications - 15th International Conference, ADMA 2019, Dalian, China, November 21-23, 2019, Proceedings*. 406–420.
- [111] Zhihui Li, Feiping Nie, Xiaoju Chang, Zhigang Ma, and Yi Yang. 2018. Balanced Clustering via Exclusive Lasso: A Pragmatic Approach. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*. 3596–3603.
- [112] Zhihui Li, Feiping Nie, Xiaoju Chang, Liqiang Nie, Huaxiang Zhang, and Yi Yang. 2018. Rank-Constrained Spectral Clustering With Flexible Embedding. *IEEE Trans. Neural Networks Learn. Syst.* 29, 12 (2018), 6073–6082.
- [113] Zhihui Li, Feiping Nie, Xiaoju Chang, and Yi Yang. 2017. Beyond Trace Ratio: Weighted Harmonic Mean of Trace Ratios for Multiclass Discriminant Analysis. *IEEE Trans. Knowl. Data Eng.* 29, 10 (2017), 2100–2110.
- [114] Zhihui Li, Feiping Nie, Xiaoju Chang, Yi Yang, Chengqi Zhang, and Nicu Sebe. 2018. Dynamic Affinity Graph Construction for Spectral Clustering Using Multiple Features. *IEEE Trans. Neural Networks Learn. Syst.* 29, 12 (2018), 6323–6332.
- [115] Zhihui Li, Lina Yao, Xiaoju Chang, Kun Zhan, Jiande Sun, and Huaxiang Zhang. 2019. Zero-shot event detection via event-adaptive concept relevance mining. *Pattern Recognit.* 88 (2019), 595–603.
- [116] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*. Springer, 740–755.
- [117] Chong Liu, Xiaoju Chang, and Yi-Dong Shen. 2020. Unity Style Transfer for Person Re-Identification. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*. 6886–6895.
- [118] Huan Liu, Qinghua Zheng, Minnan Luo, Xiaoju Chang, Caixia Yan, and Lina Yao. 2020. Memory transformation networks for weakly supervised visual classification. *Knowl. Based Syst.* 210 (2020), 106432.
- [119] Huan Liu, Qinghua Zheng, Minnan Luo, Dingwen Zhang, Xiaoju Chang, and Cheng Deng. 2017. How Unlabeled Web Videos Help Complex Event Detection?. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*. 4040–4046.
- [120] Wenhe Liu, Xiaoju Chang, Ling Chen, Dinh Hung, Xiaoju Chang, Yi Yang, and Alexander G. Hauptmann. 2020. Pair-based Uncertainty and Diversity Promoting Early Active Learning for Person Re-identification. *ACM Trans. Intell. Syst. Technol.* 11, 2 (2020), 21:1–21:15.
- [121] Wenhe Liu, Xiaoju Chang, Ling Chen, and Yi Yang. 2017. Early Active Learning with Pairwise Constraint for Person Re-identification. In *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD 2017, Skopje, Macedonia, September 18-22, 2017, Proceedings, Part I*. 103–118.
- [122] Wenhe Liu, Xiaoju Chang, Ling Chen, and Yi Yang. 2018. Semi-Supervised Bayesian Attribute Learning for Person Re-Identification. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*. 7162–7169.
- [123] Wenhe Liu, Xiaoju Chang, Yan Yan, Yi Yang, and Alexander G. Hauptmann. 2018. Few-Shot Text and Image Classification via Analogical Transfer Learning. *ACM Trans. Intell. Syst. Technol.* 9, 6 (2018), 71:1–71:20.
- [124] Wenhe Liu, Guoliang Kang, Po-Yao Huang, Xiaoju Chang, Lijun Yu, Yijun Qian, Junwei Liang, Liangke Gui, Jing Wen, Peng Chen, and Alexander G. Hauptmann. 2020. Argus: Efficient Activity Detection System for Extended Video Analysis. In *IEEE Winter Applications of Computer Vision Workshops, WACV Workshops 2020, Snowmass Village, CO, USA, March 1-5, 2020*. 126–133.
- [125] Minnan Luo, Xiaoju Chang, Zhihui Li, Liqiang Nie, Alexander G. Hauptmann, and Qinghua Zheng. 2017. Simple to complex cross-modal learning to rank. *Comput. Vis. Image Underst.* 163 (2017), 67–77.
- [126] Minnan Luo, Xiaoju Chang, Liqiang Nie, Yi Yang, Alexander G. Hauptmann, and Qinghua Zheng. 2018. An Adaptive Semisupervised Feature Analysis for Video Semantic Recognition. *IEEE Trans. Cybern.* 48, 2 (2018), 648–660.
- [127] Minnan Luo, Feiping Nie, Xiaoju Chang, Yi Yang, Alexander G. Hauptmann, and Qinghua Zheng. 2016. Avoiding Optimal Mean Robust PCA/2DPCA with Non-greedy  $\ell_1$ -Norm Maximization. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI 2016, New York, NY, USA, 9-15 July 2016*. 1802–1808.
- [128] Minnan Luo, Feiping Nie, Xiaoju Chang, Yi Yang, Alexander G. Hauptmann, and Qinghua Zheng. 2017. Avoiding Optimal Mean  $\ell_2, \ell_1$ -Norm Maximization-Based Robust PCA for Reconstruction. *Neural Comput.* 29, 4 (2017), 1124–1150.
- [129] Minnan Luo, Feiping Nie, Xiaoju Chang, Yi Yang, Alexander G. Hauptmann, and Qinghua Zheng. 2018. Adaptive Unsupervised Feature Selection With Structure Regularization. *IEEE Trans. Neural Networks Learn. Syst.* 29, 4 (2018), 944–956.
- [130] Minnan Luo, Caixia Yan, Qinghua Zheng, Xiaoju Chang, Ling Chen, and Feiping Nie. 2019. Discrete Multi-Graph Clustering. *IEEE Trans. Image Process.* 28, 9 (2019), 4701–4712.
- [131] Minnan Luo, Lingling Zhang, Feiping Nie, Xiaoju Chang, Buyue Qian, and Qinghua Zheng. 2017. Adaptive Semi-Supervised Learning with Discriminative Least Squares Regression. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*. 2421–2427.
- [132] Minh-Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. Effective Approaches to Attention-based Neural Machine Translation. In *Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational

- Linguistics, Lisbon, Portugal, 1412–1421. <http://aclweb.org/anthology/D15-1166>
- [133] Thang Luong, Hieu Pham, and Christopher D Manning. 2015. Effective Approaches to Attention-based Neural Machine Translation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. 1412–1421.
- [134] Zhigang Ma, Xiaojun Chang, Zhongwen Xu, Nicu Sebe, and Alexander G. Hauptmann. 2018. Joint Attributes and Event Analysis for Multimedia Event Detection. *IEEE Trans. Neural Networks Learn. Syst.* 29, 7 (2018), 2921–2930.
- [135] Zhigang Ma, Xiaojun Chang, Yi Yang, Nicu Sebe, and Alexander G. Hauptmann. 2017. The Many Shades of Negativity. *IEEE Trans. Multimed.* 19, 7 (2017), 1558–1568.
- [136] Pascal Mettes, Dennis C. Koelma, and Cees G.M. Snoek. 2016. The ImageNet Shuffle: Reorganized Pre-training for Video Event Detection. In *Proceedings of the 2016 ACM on International Conference on Multimedia Retrieval (ICMR '16)*. ACM, New York, NY, USA, 175–182. <https://doi.org/10.1145/2911996.2912036>
- [137] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient Estimation of Word Representations in Vector Space. *CoRR abs/1301.3781* (2013). [arXiv:1301.3781](http://arxiv.org/abs/1301.3781) <http://arxiv.org/abs/1301.3781>
- [138] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*. 3111–3119.
- [139] Niluthpol Chowdhury Mithun, Rameswar Panda, Evangelos E Papalexakis, and Amit K Roy-Chowdhury. 2018. Webly Supervised Joint Embedding for Cross-Modal Image-Text Retrieval. In *2018 ACM Multimedia Conference on Multimedia Conference*. ACM, 1856–1864.
- [140] Hyeonseob Nam, Jung-Woo Ha, and Jeonghee Kim. 2017. Dual Attention Networks for Multimodal Reasoning and Matching. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*. IEEE, 2156–2164.
- [141] Deepak Ranjan Nayak, Ratnakar Dash, Xiaojun Chang, Banshidhar Majhi, and Sambit Bakshi. 2020. Automated Diagnosis of Pathological Brain Using Fast Curvelet Entropy Features. *IEEE Trans. Sustain. Comput.* 5, 3 (2020), 416–427.
- [142] Liqiang Nie, Luming Zhang, Lei Meng, Xueming Song, Xiaojun Chang, and Xuelong Li. 2017. Modeling Disease Progression via Multisource Multitask Learners: A Case Study With Alzheimer’s Disease. *IEEE Trans. Neural Networks Learn. Syst.* 28, 7 (2017), 1508–1519.
- [143] Liqiang Nie, Luming Zhang, Yan Yan, Xiaojun Chang, Maofu Liu, and Ling Shaoling. 2017. Multiview Physician-Specific Attributes Fusion for Health Seeking. *IEEE Trans. Cybern.* 47, 11 (2017), 3680–3691.
- [144] Ram Prasad Padhy, Xiaojun Chang, Suman Kumar Choudhury, Pankaj Kumar Sa, and Sambit Bakshi. 2019. Multi-stage cascaded deconvolution for depth map and surface normal prediction from single image. *Pattern Recognit. Lett.* 127 (2019), 165–173.
- [145] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*. Association for Computational Linguistics, 311–318.
- [146] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. GloVe: Global Vectors for Word Representation. In *Empirical Methods in Natural Language Processing (EMNLP)*. 1532–1543. <http://www.aclweb.org/anthology/D14-1162>
- [147] Florent Perronnin and Christopher Dance. 2007. Fisher kernels on visual vocabularies for image categorization. In *2007 IEEE conference on computer vision and pattern recognition*. IEEE, 1–8.
- [148] Alberto Poncelas, Dimitar Shterionov, Andy Way, Gideon Maillette de Buy Wenniger, and Peyman Passban. 2018. Investigating Backtranslation in Neural Machine Translation. *arXiv preprint arXiv:1804.06189* (2018).
- [149] Janarthanan Rajendran, Mitesh M Khapra, Sarath Chandar, and Balaraman Ravindran. 2016. Bridge Correlational Neural Networks for Multilingual Multimodal Representation Learning. In *Proceedings of NAACL-HLT*. 171–181.
- [150] Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Xiaojiang Chen, and Xin Wang. 2020. A Comprehensive Survey of Neural Architecture Search: Challenges and Solutions. *CoRR abs/2006.02903* (2020).
- [151] Pengzhen Ren, Yun Xiao, Xiaojun Chang, Mahesh Prakash, Feiping Nie, Xin Wang, and Xiaojiang Chen. 2020. Structured Optimal Graph-Based Clustering With Flexible Embedding. *IEEE Trans. Neural Networks Learn. Syst.* 31, 10 (2020), 3801–3813.
- [152] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*. 91–99.
- [153] Imad Rida, Sambit Bakshi, Xiaojun Chang, and Hugo Proença. 2019. Forensic Shoe-print Identification: A Brief Survey. *CoRR abs/1901.01431* (2019).
- [154] Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Edinburgh neural machine translation systems for wmt 16. *arXiv preprint arXiv:1606.02891* (2016).
- [155] Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving Neural Machine Translation Models with Monolingual Data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Vol. 1. 86–96.
- [156] Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural Machine Translation of Rare Words with Subword Units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Vol. 1. 1715–1725.
- [157] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [158] Amit Kumar Singh, Zhihan Lv, Huimin Lu, and Xiaojun Chang. 2020. Guest editorial: Recent trends in multimedia data-hiding: a reliable mean for secure communications. *J. Ambient Intell. Humaniz. Comput.* 11, 5 (2020), 1795–1797.
- [159] Chen Sun, Sanketh Shetty, Rahul Sukthankar, and Ram Nevatia. 2015. Temporal localization of fine-grained actions in videos by domain transfer from web images. In *Proceedings of the 23rd ACM international conference on Multimedia*. ACM, 371–380.
- [160] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to Sequence Learning with Neural Networks. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*. 3104–3112.
- [161] Bart Thomee, David A. Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. 2016. YFCC100M: The New Data in Multimedia Research. *Commun. ACM* 59, 2 (2016), 64–73. <http://cacm.acm.org/magazines/2016/2/197425-yfcc100m/fulltext>
- [162] Yao-Hung Hubert Tsai, Liang-Kang Huang, and Ruslan Salakhutdinov. 2017. Learning robust visual-semantic embeddings. In *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE, 3591–3600.
- [163] Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew W. Senior, and Koray Kavukcuoglu. 2016. WaveNet: A Generative Model for Raw Audio. In *The 9th ISCA Speech Synthesis Workshop, Sunnyvale, CA, USA, 13-15 September 2016*. 125.
- [164] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*. 5998–6008.
- [165] Ivan Vendrov, Ryan Kiros, Sanja Fidler, and Raquel Urtasun. 2015. Order-embeddings of images and language. *arXiv preprint arXiv:1511.06361* (2015).
- [166] Ivan Vendrov, Ryan Kiros, Sanja Fidler, and Raquel Urtasun. 2015. Order-Embeddings of Images and Language. *CoRR abs/1511.06361* (2015).
- [167] Bokun Wang, Yang Yang, Xing Xu, Alan Hanjalic, and Heng Tao Shen. 2017. Adversarial Cross-Modal Retrieval. In *ACM MM*.
- [168] Fei Wang, Lei Zhu, Cheng Liang, Jingjing Li, Xiaojun Chang, and Ke Lu. 2020. Robust optimal graph clustering. *Neurocomputing* 378 (2020), 153–165.
- [169] Hanno Wang, Xiaojun Chang, Lei Shi, Yi Yang, and Yi-Dong Shen. 2018. Uncertainty Sampling for Action Recognition via Maximizing Expected Average Precision. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*. 964–970.
- [170] Kaiye Wang, Qiyue Yin, Wei Wang, Shu Wu, and Liang Wang. 2016. A comprehensive survey on cross-modal retrieval. *arXiv preprint arXiv:1607.06215* (2016).
- [171] Liwei Wang, Yin Li, Jing Huang, and Svetlana Lazebnik. 2018. Learning two-branch neural networks for image-text matching tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2018).
- [172] Liwei Wang, Yin Li, Jing Huang, and Svetlana Lazebnik. 2019. Learning two-branch neural networks for image-text matching tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41, 2 (2019), 394–407.
- [173] Liwei Wang, Yin Li, and Svetlana Lazebnik. 2016. Learning Deep Structure-Preserving Image-Text Embeddings. In *Computer Vision and Pattern Recognition (CVPR), 2016 IEEE Conference on*. IEEE, 5005–5013.
- [174] Rong Wang, Feiping Nie, Richang Hong, Xiaojun Chang, Xiaojun Yang, and Weizhong Yu. 2017. Fast and Orthogonal Locality Preserving Projections for Dimensionality Reduction. *IEEE Trans. Image Process.* 26, 10 (2017), 5019–5030.
- [175] Sen Wang, Xiaojun Chang, Xue Li, Guodong Long, Lina Yao, and Quan Z. Sheng. 2016. Diagnosis Code Assignment Using Sparsity-Based Disease Correlation Embedding. *IEEE Trans. Knowl. Data Eng.* 28, 12 (2016), 3191–3202.
- [176] Sen Wang, Xue Li, Xiaojun Chang, Lina Yao, Quan Z. Sheng, and Guodong Long. 2017. Learning Multiple Diagnosis Codes for ICU Patients with Local Disease Correlation Mining. *ACM Trans. Knowl. Discov. Data* 11, 3 (2017), 31:1–31:21.
- [177] Sen Wang, Feiping Nie, Xiaojun Chang, Xue Li, Quan Z. Sheng, and Lina Yao. 2016. Uncovering Locally Discriminative Structure for Feature Analysis. In *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD 2016, Riva del Garda, Italy, September 19-23, 2016, Proceedings, Part I*. 281–295.
- [178] Sen Wang, Feiping Nie, Xiaojun Chang, Lina Yao, Xue Li, and Quan Z. Sheng. 2015. Unsupervised Feature Analysis with Class Margin Optimization. In *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD 2015, Porto, Portugal, September 7-11, 2015, Proceedings, Part I*. 383–398.
- [179] Weitao Wang, Ruyang Liu, Meng Wang, Sen Wang, Xiaojun Chang, and Yang Chen. 2020. Memory-Based Network for Scene Graph with Unbalanced Relations. In *MM '20: The 28th ACM International Conference on Multimedia, Virtual*

- Event / Seattle, WA, USA, October 12-16, 2020. 2400–2408.
- [180] Man Wu, Shirui Pan, Chuan Zhou, Xiaojun Chang, and Xingquan Zhu. 2020. Unsupervised Domain Adaptive Graph Convolutional Networks. In *WWW '20: The Web Conference 2020, Taipei, Taiwan, April 20-24, 2020*. 1457–1467.
- [181] Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, Xiaojun Chang, and Chengqi Zhang. 2020. Connecting the Dots: Multivariate Time Series Forecasting with Graph Neural Networks. In *KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020*. 753–763.
- [182] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*. 2048–2057.
- [183] Shicheng Xu, Huan Li, Xiaojun Chang, Shouo-I Yu, Xingzhong Du, Xuanchong Li, Lu Jiang, Zexi Mao, Zhen-Zhong Lan, Susanne Burger, and Alexander G. Hauptmann. 2015. Incremental Multimodal Query Construction for Video Search. In *Proceedings of the 5th ACM on International Conference on Multimedia Retrieval, Shanghai, China, June 23-26, 2015*. 675–678.
- [184] Xing Xu, Li He, Huimin Lu, Lianli Gao, and Yanli Ji. [n. d.]. Deep adversarial metric learning for cross-modal retrieval. *World Wide Web* ([n. d.]), 1–16.
- [185] Xiaowei Xue, Feiping Nie, Sen Wang, Xiaojun Chang, Bela Stantic, and Min Yao. 2017. Multi-View Correlated Feature Learning by Uncovering Shared Component. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*. 2810–2816.
- [186] Caixia Yan, Xiaojun Chang, Minnan Luo, Qinghua Zheng, Xiaoqin Zhang, Zhihui Li, and Feiping Nie. 2020. Self-Weighted Robust LDA for Multiclass Classification with Edge Classes. *CoRR abs/2009.12362* (2020).
- [187] Caixia Yan, Qinghua Zheng, Xiaojun Chang, Minnan Luo, Chung-Hsing Yeh, and Alexander G. Hauptmann. 2020. Semantics-Preserving Graph Propagation for Zero-Shot Object Detection. *IEEE Trans. Image Process.* 29 (2020), 8163–8176.
- [188] Fei Yan and Krystian Mikolajczyk. 2015. Deep correlation for matching images and text. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3441–3450.
- [189] Yi Yang, Zhigang Ma, Feiping Nie, Xiaojun Chang, and Alexander G. Hauptmann. 2015. Multi-Class Active Learning by Uncertainty Sampling with Diversity Maximization. *Int. J. Comput. Vis.* 113, 2 (2015), 113–127.
- [190] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics* 2 (2014), 67–78.
- [191] En Yu, Wenhe Liu, Guoliang Kang, Xiaojun Chang, Jiande Sun, and Alexander G. Hauptmann. 2019. Inf@TRECVID 2019: Instance Search Task. In *2019 TREC Video Retrieval Evaluation, TRECVID 2019, Gaithersburg, MD, USA, November 12-13, 2019*.
- [192] En Yu, Jiande Sun, Jing Li, Xiaojun Chang, Xian-Hua Han, and Alexander G. Hauptmann. 2019. Adaptive Semi-Supervised Feature Selection for Cross-Modal Retrieval. *IEEE Trans. Multimed.* 21, 5 (2019), 1276–1288.
- [193] En Yu, Jiande Sun, Li Wang, Xiaojun Chang, Huaxiang Zhang, and Alexander G. Hauptmann. 2019. Cross-Modal Transfer Hashing Based on Coherent Projection. In *IEEE International Conference on Multimedia & Expo Workshops, ICME Workshops 2019, Shanghai, China, July 8-12, 2019*. 477–482.
- [194] Shouo-I Yu, Lu Jiang, Zhongwen Xu, Zhenzhong Lan, Shicheng Xu, Xiaojun Chang, Xuanchong Li, Zexi Mao, Chuang Gan, Yajie Miao, Xingzhong Du, Yang Cai, Lara J. Martin, Nikolas Wolfe, Anurag Kumar, Huan Li, Ming Lin, Zhigang Ma, Yi Yang, Deyu Meng, Shiguang Shan, Pinar Duygulu Sahin, Susanne Burger, Florian Metz, Rita Singh, Bhiksha Raj, Teruko Mitamura, Richard M. Stern, and Alexander G. Hauptmann. 2015. CMU Informedia@TRECVID 2015: MED/SIN/LNK/SED. In *2015 TREC Video Retrieval Evaluation, TRECVID 2015, Gaithersburg, MD, USA, November 16-18, 2015*. <https://www.nlpir.nist.gov/projects/tvpubs/tv15.papers/cmu.pdf>
- [195] Zhen Yu, Jennifer Nguyen, Xiaojun Chang, John Kelly, Catriona McLean, Lei Zhang, Victoria Mar, and Zongyuan Ge. 2020. Melanoma Diagnosis with Spatio-Temporal Feature Learning on Sequential Dermoscopic Images. *CoRR abs/2006.10950* (2020).
- [196] Di Yuan, Xiaojun Chang, and Zhenyu He. 2020. Accurate Bounding-box Regression with Distance-LoU Loss for Visual Tracking. *CoRR abs/2007.01864* (2020).
- [197] Di Yuan, Xiaojun Chang, Po-Yao Huang, Qiao Liu, and Zhenyu He. 2021. Self-Supervised Deep Correlation Tracking. *IEEE Trans. Image Process.* 30 (2021), 976–985. <https://doi.org/10.1109/TIP.2020.3037518>
- [198] Kun Zhan, Xiaojun Chang, Junpeng Guan, Ling Chen, Zhigang Ma, and Yi Yang. 2019. Adaptive Structure Discovery for Multimedia Analysis Using Multiple Features. *IEEE Trans. Cybern.* 49, 5 (2019), 1826–1834.
- [199] Dingwen Zhang, Junwei Han, Lu Jiang, Senmao Ye, and Xiaojun Chang. 2017. Revealing Event Saliency in Unconstrained Video Collection. *IEEE Trans. Image Process.* 26, 4 (2017), 1746–1758.
- [200] Dalin Zhang, Lina Yao, Kaixuan Chen, Sen Wang, Xiaojun Chang, and Yunhao Liu. 2020. Making Sense of Spatio-Temporal Preserving Representations for EEG-Based Human Intention Recognition. *IEEE Trans. Cybern.* 50, 7 (2020), 3033–3044.
- [201] Jiaqi Zhang, Meng Wang, Qinchu Li, Sen Wang, Xiaojun Chang, and Beilun Wang. 2020. Quadratic Sparse Gaussian Graphical Model Estimation Method for Massive Variables. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*. 2964–2972.
- [202] Lingling Zhang, Xiaojun Chang, Jun Liu, Minnan Luo, Mahesh Prakash, and Alexander G. Hauptmann. 2020. Few-shot activity recognition with cross-modal memory network. *Pattern Recognit.* 108 (2020), 107348.
- [203] Lingling Zhang, Xiaojun Chang, Jun Liu, Minnan Luo, Sen Wang, Zongyuan Ge, and Alexander G. Hauptmann. 2020. ZSTAD: Zero-Shot Temporal Activity Detection. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*. 876–885.
- [204] Lingling Zhang, Jun Liu, Minnan Luo, Xiaojun Chang, and Qinghua Zheng. 2018. Deep Semisupervised Zero-Shot Learning with Maximum Mean Discrepancy. *Neural Comput.* 30, 5 (2018).
- [205] Lingling Zhang, Jun Liu, Minnan Luo, Xiaojun Chang, Qinghua Zheng, and Alexander G. Hauptmann. 2019. Scheduled sampling for one-shot learning via matching network. *Pattern Recognit.* 96 (2019).
- [206] Lingling Zhang, Minnan Luo, Jun Liu, Xiaojun Chang, Yi Yang, and Alexander G. Hauptmann. 2020. Deep Top- $k$  Ranking for Image-Sentence Matching. *IEEE Trans. Multimed.* 22, 3 (2020), 775–785.
- [207] Miao Zhang, Huiqi Li, Shirui Pan, Xiaojun Chang, Zongyuan Ge, and Steven W. Su. 2020. Differentiable Neural Architecture Search in Equivalent Space with Exploration Enhancement. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- [208] Miao Zhang, Huiqi Li, Shirui Pan, Xiaojun Chang, and Steven W. Su. 2020. Overcoming Multi-Model Forgetting in One-Shot NAS With Diversity Maximization. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*. 7806–7815.
- [209] Zhicheng Zhao, Xuanchong Li, Xingzhong Du, Qi Chen, Yanyun Zhao, Fei Su, Xiaojun Chang, and Alexander G. Hauptmann. 2018. A unified framework with a benchmark dataset for surveillance event detection. *Neurocomputing* 278 (2018), 62–74.
- [210] Zhedong Zheng, Liang Zheng, Michael Garrett, Yi Yang, and Yi-Dong Shen. 2017. Dual-Path Convolutional Image-Text Embedding. *CoRR abs/1711.05535* (2017). <http://arxiv.org/abs/1711.05535>
- [211] Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. 2014. Learning Deep Features for Scene Recognition using Places Database. In *Advances in Neural Information Processing Systems 27, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger (Eds.), Curran Associates, Inc., 487–495*. <http://papers.nips.cc/paper/5349-learning-deep-features-for-scene-recognition-using-places-database.pdf>
- [212] Runwu Zhou, Xiaojun Chang, Lei Shi, Yi-Dong Shen, Yi Yang, and Feiping Nie. 2020. Person Reidentification via Multi-Feature Fusion With Adaptive Graph Learning. *IEEE Trans. Neural Networks Learn. Syst.* 31, 5 (2020), 1592–1601.
- [213] Fengda Zhu, Xiaojun Chang, Runhao Zeng, and Mingkui Tan. 2019. Continual Reinforcement Learning with Diversity Exploration and Adversarial Self-Correction. *CoRR abs/1906.09205* (2019).
- [214] Fengda Zhu, Yi Zhu, Xiaojun Chang, and Xiaodan Liang. 2020. Vision-Language Navigation With Self-Supervised Auxiliary Reasoning Tasks. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*. 10009–10019.
- [215] Lei Zhu, Zi Huang, Xiaojun Chang, Jingkuan Song, and Heng Tao Shen. 2017. Exploring Consistent Preferences: Discrete Hashing with Pair-Exemplar for Scalable Landmark Search. In *Proceedings of the 2017 ACM on Multimedia Conference, MM 2017, Mountain View, CA, USA, October 23-27, 2017*. 726–734.
- [216] Yi Zhu, Fengda Zhu, Zhaohuan Zhan, Bingqian Lin, Jianbin Jiao, Xiaojun Chang, and Xiaodan Liang. 2020. Vision-Dialog Navigation by Exploring Cross-Modal Memory. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*. 10727–10736.