

---

# Multi-source Multi-modal Domain Adaptation for Visual-textual Sentiment Classification

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Learning from multiple modalities has recently been paid increasing attention  
2 in sentiment analysis tasks because of its ability to capture the complementary  
3 representation of the intrinsic multi-modal world. Recent deep learning-based  
4 multi-modal sentiment analysis methods trained on large-scale labeled data cannot  
5 guarantee good generalization to another target domain, because of the presence of  
6 domain shift. Multi-modal domain adaptation (MMDA) aims to address this issue  
7 by learning a transferable model with specific alignment across domains. However,  
8 existing MMDA methods only focus on the single-source scenario with only one  
9 labeled source domain. When labeled data is collected practically from multiple  
10 sources with different distributions, naive application of these single-source MMDA  
11 methods would fail without considering the domain shift among different sources.  
12 In this paper, we propose to study multi-source MMDA for visual-textual sentiment  
13 classification and design a novel multi-source multi-modal contrastive adversarial  
14 network, termed M2CAN, to learn domain-invariant multi-modal representations.  
15 Specifically, the designed M2CAN jointly optimizes three different alignment  
16 strategies: cross-modal contrastive alignment within each domain, cross-domain  
17 contrastive alignment for each modality, and cross-domain adversarial alignment  
18 on the fused multi-modal representations. After such alignments, different source  
19 and target domains are mapped into a shared multi-modal representation space. We  
20 conduct extensive experiments on a benchmark dataset with three domains and the  
21 results demonstrate that the proposed M2CAN significantly outperforms state-of-  
22 the-art domain adaptation approaches for visual-textual sentiment classification.  
23 Our source code will be released.

## 24 1 Introduction

25 Customers have become used to sharing their experiences and opinions of the products and service  
26 they purchase by posting reviews or comments on social networks [1, 2, 3, 4, 5, 6]. Sentiment analysis  
27 of the large-scale user-generated multimedia data plays a vitally important role in both customers’  
28 product selection and enterprises’ product improvement. On the one hand, it can influence customers’  
29 decision-making when selecting what they want. For example, if the feedback from other customers  
30 is dominated by negative comments, it is highly probable that the current customers change their  
31 attitudes to another brand. On the other hand, it can help enterprises to analyze the drawbacks revealed  
32 by customers and correspondingly improve the quality of their products and services [1, 7]. Although  
33 text is one direct and popular modality to express customers’ opinions [2], sentiment analysis solely  
34 from text may not well reflect the customers’ actual feelings. For example, if we see a comment like  
35 “what a good restaurant!”, we may conclude that the customer is satisfied with the dining; but if there  
36 is an also an affiliated image showing a dirty and disorderly environment, we can infer that the text is  
37 actually sarcasm and that the customer is upset about the dining environment. Therefore, sentiment

38 analysis from multiple modalities, such as image and text, has attracted increasing research attention  
39 with the help of easy photographing on mobile devices.

40 Recently, deep neural networks (DNNs) have achieved the state-of-the-art performances on visual-  
41 textual sentiment classification by effectively exploring the abundant and complementary content  
42 knowledge from different modalities [8, 9, 10, 11, 12, 13]. To train a DNN well, large-scale  
43 annotations are often required; however, these are not always available, since labeling multi-modal  
44 data is time-consuming and even difficult. One may consider transferring the trained DNN on a labeled  
45 source domain to the unlabeled target domain as an alternate solution. Obviously, direct transfer  
46 cannot guarantee good generalization and often results in significant performance decay [14, 15, 16],  
47 because of the presence of domain shift [17], *i.e.* the distributions of observed multi-modal data and  
48 sentiment are different between the source and target domains. Aiming at minimizing the domain  
49 gap, domain adaptation (DA) [18, 19, 20, 21, 22] tries to learn a model on the labeled source domain  
50 that can generalize well to the target domain through specific alignment across domains, such as  
51 discrepancy-based, adversarial, and self-supervision-based methods.

52 Current DA methods for visual-textual sentiment classification and other multi-modal learning  
53 tasks only focus on the single-source unsupervised setting [23, 24, 25, 26, 27], by assuming that  
54 the labeled source data is collected from the same distribution. However, in practice, it is more  
55 practical that the labeled multi-modal data comes from different source distributions [19, 21]. For  
56 example, user-generated reviews can be collected from Yelp, Twitter, and Amazon. We can naively  
57 combine different sources into one source and directly apply existing single-source DA algorithms.  
58 However, because of the neglect of mis-alignment across different sources, such methods may lead to  
59 sub-optimal results [21] (see the comparison between single-best and source-combined MM-SADA  
60 in Table 1). Therefore, effective multi-source domain adaptation (MSDA) techniques [19, 21] are  
61 required to sufficiently leverage the complementary information from different sources.

62 Recently, some deep MSDA methods have been proposed. Based on different alignment strategies,  
63 Zhao et al. classified them into two categories [21], *i.e.* latent space transformation [28, 29, 30, 31,  
64 32, 33, 34, 35, 36, 16, 37, 38, 39] and intermediate domain generation [40, 41, 42, 43, 44]. All these  
65 MSDA methods only consider a single modality, such as text or image. When extending them to  
66 a multi-modal setting, they usually fail since they cannot deal well with the heterogeneity gap, *i.e.*  
67 the semantic difference between data in different modalities (*e.g.* heterogeneity of the feature space  
68 of each modality and data content) [23]. Therefore, ineffectively aligning feature representations  
69 and mining cross-modal information may result in interference among different modalities, leading  
70 classification models to fail to capture accurate and stable sentiment-related patterns.

71 In this paper, we generalize the single-source MMDA and single-modal MSDA problems to multi-  
72 source multi-modal domain adaptation (MS-MMDA) problem, and design a novel multi-source multi-  
73 modal contrastive adversarial network, termed M2CAN, for visual-textual sentiment classification.  
74 First, we use a pair of pre-trained image and text encoders in order to project images and texts from  
75 different domains into a continuous latent feature space. Second, we perform different alignments to  
76 learn domain-invariant multi-modal representations, including (1) cross-modal contrastive alignment  
77 on the transformed lower-dimensional representations obtained by a non-linear transformation layer  
78 within each domain, (2) cross-domain contrastive alignment on the original representations for  
79 each modality, and (3) cross-domain adversarial alignment on the fused multi-modal representations  
80 obtained by multi-modal low-rank bi-linear pooling. Finally, we train a transferable task sentiment  
81 classifier based on the aligned multi-modal feature representations and corresponding source labels.  
82 Extensive experiments are conducted on a combined dataset consisting of three domains, *i.e.* Yelp [12],  
83 Twitter [45], and MVSA [46]. The results demonstrate that M2CAN significantly outperforms the  
84 state-of-the-art DA methods for visual-textual sentiment classification.

85 In summary, the contributions of this paper are threefold: (1) We propose to study a novel and  
86 practical DA setting, *i.e.* multi-source multi-modal domain adaptation (MS-MMDA), for visual-  
87 textual sentiment classification. To the best of our knowledge, this is the first work that investigates  
88 MMDA with multiple sources. (2) We propose a novel MS-MMDA method, termed M2CAN,  
89 by contrastive and adversarial learning. Through both cross-modal alignment and cross-domain  
90 alignment, M2CAN can learn domain invariant multi-modal representations and thus minimizes  
91 the domain gap among multiple sources and the target. (3) We conduct extensive experiments on  
92 a benchmark dataset with three different domains. As compared to the best baseline, the proposed  
93 M2CAN achieves 2.7% performance gains on the average classification accuracy.

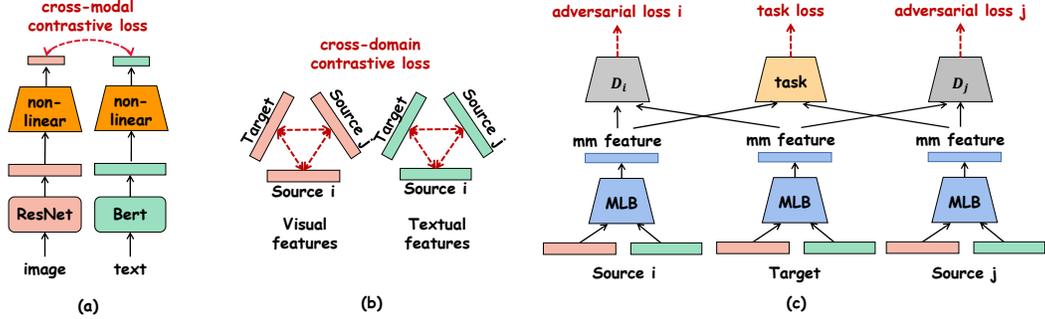


Figure 1: Illustration of the proposed M2CAN framework: (a) image-text feature encoding and cross-modal contrastive alignment, (b) cross-domain contrastive alignment, and (c) cross-domain adversarial alignment and task classifier learning. All images and texts are encoded with encoders (ResNet50 [47] and BERT [48]) to a latent continuous feature space. Three different alignments are then performed to learn domain-invariant multi-modal representations, including cross-modal transformation, cross-domain contrastive alignment on the original representations for each modality, and cross-domain adversarial alignment on the fused multi-modal representations obtained by multi-modal low-rank bi-linear pooling (MLB). A transferable task sentiment classifier is finally trained based on the aligned multi-modal (mm) feature representations and corresponding source labels.

## 94 2 Multi-source Multi-modal Domain Adaptation Network

95 We consider the multi-source domain adaptation setup for visual-textual sentiment classification,  
 96 under the *covariate shift* assumption [18]. Assume access to  $K$  source domains  $\{\mathcal{S}_i\}_{i=1}^K$  with  
 97 labeled training data and a target domain  $\mathcal{T}$  with unlabeled training data consisting of two modalities,  
 98 *i.e.* image and text. Each domain  $\mathcal{S}_i$  contains a set of examples drawn from a joint distribution  
 99  $p^{(\mathcal{S}_i)}(x_{text}, x_{image}, y)$  on the input space  $\mathcal{X}_{text} \times \mathcal{X}_{image}$  and the output space  $\mathcal{Y}$ , and we seek to  
 100 learn a sentiment classifier  $f : \mathcal{X}_{text} \times \mathcal{X}_{image} \rightarrow \mathcal{Y}$  that is transferable to a target domain  $\mathcal{T}$ , where  
 101 only unlabeled data is available. In this section, we give an overview of M2CAN, present each  
 102 component of M2CAN in detail, and finally introduce the joint learning process.

### 103 2.1 Overview

104 The proposed M2CAN bridges the domain gap by performing both contrastive and adversarial  
 105 alignments among the source and target domains. The framework is shown in Figure 1. In addition to  
 106 the pre-trained encoders to encode texts and images from different domains into a semantic-preserving  
 107 latent continuous feature space and the task classifier to train the final sentiment classification model  
 108 based on the aligned multi-modal features, it consists of three primary alignment components:

109 *Cross-modal contrastive alignment (CMCA)*: Align the encoded lower-dimensional representations  
 110 between different modalities within each domain. The visual and textual representations are projected  
 111 into a lower-dimensional space with a non-linear transformation layer to extract data transformation-  
 112 invariant features. A contrastive loss is employed to align visual and textual representations by  
 113 minimizing the spatial distance between related image and text and maximizing the distance between  
 114 unrelated pairs.

115 *Cross-domain contrastive alignment (CDCA)*: Align the encoded original representations between  
 116 different domains for each modality. Considering that domain gap exists in each modality, the  
 117 discrepancy between different domains is decreased for each modality through contrastive learning.  
 118 Due to the fact that there are no negative pairs of samples, applying the same mechanism with  
 119 CMCA can lead to mode collapse of non-linear layer, *i.e.* the non-linear layer might tend to project  
 120 dissimilar higher dimensional features into similar lower dimensional features. Therefore, the CDCA  
 121 is constructed on the original feature space.

122 *Cross-domain adversarial alignment (CDAA)*: Align the fused multi-modal representations between  
 123 different domains. A fused multi-modal feature space  $\mathcal{X}_{mm}$  is created by using a bi-linear pooling

124 layer, which learns a semantic-preserving and semantic-relevant projection  $\mathcal{X}_{text} \times \mathcal{X}_{image} \rightarrow \mathcal{X}_{mm}$ .  
 125 Adversarial learning is employed to align the fused multi-modal features from different domains.

## 126 2.2 Cross-modal Contrastive Alignment

127 Simply extracting visual and textual features using separate encoders does not take the discrepancy  
 128 between features in different modalities into account. Practically, user-generated data of visual-textual  
 129 pairs might contain unrelated sentiment information. Furthermore, unaligned visual and textual  
 130 features are from different feature spaces, which might affect the sentiment classification network’s  
 131 ability to learn appropriate patterns related to sentiment. Therefore, alignment between features from  
 132 multiple modalities is necessary. For this purpose, we follow [49] and incorporate a contrastive loss  
 133 into our network. By applying data augmentation in both images and texts, we construct positive  
 134 and negative sample pairs of different modalities in each domain on a lower-dimensional space  
 135 using a non-linear transformation layer. Since the non-linear layer is trained to be invariant to data  
 136 transformation, it can remove information that may be useful for the downstream task, such as the  
 137 color and orientation of objects, or tone-related words. By leveraging the non-linear transformation,  
 138 more information can be formed and maintained in the original features [49]. Assuming we have a  
 139 batch of visual features  $I$ , and corresponding batch of textual features  $T$ , after data augmentation, the  
 140 corresponding batches of visual and textual features are  $I'$  and  $T'$ , the cross-modal contrastive loss  
 141 can be constructed as follows [50]:

$$\mathbb{I} = g(X_{img}), \mathbb{I}' = g(X'_{img}), \mathbb{T} = g(X_{txt}), \mathbb{T}' = g(X'_{txt}), \quad (1)$$

$$\mathcal{L}_{CMCA} = -\frac{1}{n} \cdot \mathbb{I}^T \cdot \log \left[ \frac{e^{\mathbb{I} \circ \mathbb{T}} + e^{\mathbb{I} \circ \mathbb{T}'} + e^{\mathbb{I}' \circ \mathbb{T}} + e^{\mathbb{I}' \circ \mathbb{T}'}}{\mathbb{1}^T \cdot (e^{\mathbb{I} \cdot \mathbb{T}^T} + e^{\mathbb{I} \cdot \mathbb{T}'^T} + e^{\mathbb{I}' \cdot \mathbb{T}^T} + e^{\mathbb{I}' \cdot \mathbb{T}'^T}) \cdot \mathbb{1}} \right], \quad (2)$$

142 where  $g : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$  is a lower-dimensional projection function,  $d$  represents the dimension of  
 143 original visual and textual feature, while  $d'$  represents the dimension after projection,  $X_{img} \in \mathbb{R}^{n \times d}$ ,  
 144  $X_{txt} \in \mathbb{R}^{n \times d}$ ,  $X'_{img} \in \mathbb{R}^{n \times d}$ , and  $X'_{txt} \in \mathbb{R}^{n \times d}$  represent a batch of original visual and textual  
 145 feature, and a batch of augmented visual and textual feature respectively,  $\circ$  represents the Hadamard  
 146 product, and  $n$  denotes the batch size. By minimizing the distance between visual and textual features  
 147 from the same sample and maximizing the distance between visual and textual feature from different  
 148 samples before and after data augmentation, our cross-modal contrastive loss is able to force the  
 149 encoders to extract closer features from semantically similar samples and farther apart features from  
 150 semantically different samples robustly, which achieves the purpose of CMCA.

## 151 2.3 Cross-domain Contrastive Alignment

152 To describe the similarity of two distributions, we introduce the Maximum Mean Discrepancy (MMD),  
 153 as described below:

$$\mathcal{D}_{\mathcal{H}}(P, Q) \triangleq \sup_{f \in \mathcal{H}} (\mathbb{E}_{\mathbf{X}^s} [f(\mathbf{X}^s)] - \mathbb{E}_{\mathbf{X}^t} [f(\mathbf{X}^t)]), \quad (3)$$

154 where  $\mathbf{X}^s$  and  $\mathbf{X}^t$  are sampled from the marginal distributions  $P(X^s)$  and  $Q(X^t)$  respectively,  $\mathcal{H}$  is  
 155 a class of function. Formally, MMD defines the difference between two distributions with their mean  
 156 representations in the reproducing kernel Hilbert space (RKHS) [51]. In practice, the squared value  
 157 of MMD is estimated with the empirical kernel mean representations:

$$\begin{aligned} \hat{\mathcal{D}}^{mmd} &= \frac{1}{n_s^2} \sum_{i=1}^{n_s} \sum_{j=1}^{n_s} k(\phi(\mathbf{x}_i^s), \phi(\mathbf{x}_j^s)) + \frac{1}{n_t^2} \sum_{i=1}^{n_t} \sum_{j=1}^{n_t} k(\phi(\mathbf{x}_i^t), \phi(\mathbf{x}_j^t)) \\ &\quad - \frac{2}{n_s n_t} \sum_{i=1}^{n_s} \sum_{j=1}^{n_t} k(\phi(\mathbf{x}_i^s), \phi(\mathbf{x}_j^t)), \end{aligned} \quad (4)$$

158 where  $\mathbf{x}^s \in \mathcal{S}$ ,  $\mathbf{x}^t \in \mathcal{T}$ ,  $n_s$  and  $n_t$  denote the batch sizes, and  $k$  denotes a kernel function. We adopt  
 159 the third term and ignore the first two terms in Eq. (4). Due to the existence of multiple modalities, we  
 160 decompose the gap between two domains into two parts, *i.e.* visual domain gap and textual domain

161 gap, and minimize the corresponding MMD:

$$\begin{aligned} \mathcal{L}_{CDCA} = & \sum_{s_1, s_2} \sum_{\mathcal{I}^{s_1}, \mathcal{I}^{s_2}, \mathcal{T}^{s_1}, \mathcal{T}^{s_2}} - \frac{2}{n_{s_1} n_{s_2}} \sum_{i=1}^{n_{s_1}} \sum_{j=1}^{n_{s_2}} k(\phi(\mathcal{I}_i^{s_1}), \phi(\mathcal{I}_j^{s_2})) \\ & - \frac{2}{n_{s_1} n_{s_2}} \sum_{i=1}^{n_{s_1}} \sum_{j=1}^{n_{s_2}} k(\phi(\mathcal{T}_i^{s_1}), \phi(\mathcal{T}_j^{s_2})), \end{aligned}$$

$$s_1 \in Dom, s_2 \in Dom \setminus s_1, Dom = \{S_1, S_2, \dots, S_K, T\}, \mathcal{I}_i^s \in \mathcal{I}^s \in X_{img}^s \cup X_{img}^{s'}, \mathcal{T}_i^s \in \mathcal{T}^s \in X_{txt}^s \cup X_{txt}^{s'}, \quad (5)$$

162 where  $X_{img}^s$  and  $X_{img}^{s'}$  are all the possible visual feature batches in domain  $s$  with and without data  
163 augmentation,  $X_{txt}^s$  and  $X_{txt}^{s'}$  are all the possible textual feature batches in domain  $s$  with and  
164 without data augmentation,  $s \in \{s_1, s_2, \dots, s_K\}$ . We choose the linear kernel as  $k$ . Therefore, the  
165 above cross-domain contrastive loss can be simplified as below:

$$\mathcal{L}_{CDCA} = \sum_{s_1, s_2} \sum_{\mathcal{I}^{s_1}, \mathcal{I}^{s_2}, \mathcal{T}^{s_1}, \mathcal{T}^{s_2}} - \frac{2}{n_{s_1} n_{s_2}} \cdot \mathbb{1}^T \cdot \mathcal{I}_{s_1} \cdot \mathcal{I}_{s_2}^T \cdot \mathbb{1} - \frac{2}{n_{s_1} n_{s_2}} \cdot \mathbb{1}^T \cdot \mathcal{T}_{s_1} \cdot \mathcal{T}_{s_2}^T \cdot \mathbb{1}. \quad (6)$$

166 Since we cannot construct negative sample pairs as in CMCA, optimizing the above function on  
167 lower-dimensional space will cause the lower-dimensional projection function  $g$  to project all visual  
168 and textual features to all-zero features, resulting in feature encoders that are not able to learn any  
169 useful pattern. Therefore, the above optimization problem is run with on the original visual and  
170 textual features.

## 171 2.4 Cross-domain Adversarial Alignment

172 To better fuse the visual and textual features and produce a multi-modal feature space that includes  
173 enough sentiment-related visual-textual information, we choose a bi-linear model [52] to fuse each  
174 pair of features from different modalities into factors related with sentiment:

$$f_i = \sum_{p=1}^N \sum_{q=1}^M w_{ipq} \mathcal{I}^p \mathcal{T}^q + b_i = \mathcal{I}^T \mathbf{W}_i \mathcal{T} + b_i, \quad (7)$$

175 where  $\mathcal{I}$  and  $\mathcal{T}$  are visual and textual features, and  $N$  and  $M$  represent the dimension of feature  $\mathcal{I}$   
176 and  $\mathcal{T}$ , respectively.  $\mathbf{W}_i \in \mathbb{R}^{N \times M}$  represents the weight matrix for output  $f_i$ , and  $b_i$  represents the  
177 bias. Assuming the dimension of output feature is  $L$ , the number of parameters of bi-linear model is  
178  $L \times (N \times M + 1)$  including bias vector  $b$ . According to a low-rank bi-linear method which is able  
179 to reduce the dimension of the weight matrix [53], the weight matrix can be decomposed into the  
180 product of two low-order matrices, which can be described as:  $W_i = U_i V_i^T$ , where  $U_i \in \mathbb{R}^{N \times d}$  and  
181  $V_i \in \mathbb{R}^{M \times d}$ . Therefore, the output feature  $f_i$  can be formalized as:

$$f_i = \mathcal{I}^T \mathbf{W}_i \mathcal{T} + b_i = \mathcal{I}^T \mathbf{U}_i \mathbf{V}_i^T \mathcal{T} + b_i = \mathbb{1}^T (\mathbf{U}_i^T \mathcal{I} \circ \mathbf{V}_i^T \mathcal{T}) + b_i, \quad (8)$$

182 where  $\mathbb{1}$  represents a column vector consisting of component 1. Still, we need two third-order  
183 tensors,  $\mathbf{U}$  and  $\mathbf{V}$ , for a feature vector  $f$ , whose elements are  $\{f_i\}$ . To reduce the order of the weight  
184 tensors by one and introduce non-linear activation function, we adopt the following bi-linear pooling  
185 function [54]:

$$\mathbf{f} = \mathbf{P}^T (\sigma(\mathbf{U}^T \mathcal{I}) \circ \sigma(\mathbf{V}^T \mathcal{T})) + h_{img}(\mathcal{I}) + h_{txt}(\mathcal{T}) + \mathbf{b}, \quad (9)$$

186 where  $\mathbf{f}$  represents our multi-modal feature,  $\sigma$  is a non-linear activate function, and  $h_x$  and  $h_y$  are  
187 shortcut mappings.

188 In order to bridge the domain gap across multiple source domains and the target domain in the  
189 fused multi-modal feature space, we construct cross-domain adversarial alignment. Specifically, we  
190 introduce a set of domain classifiers as discriminators, which are used to distinguish source features  
191 from target features for each source. By assuming that the above encoders and bi-linear pooling  
192 layer as a feature extractor, we can construct an adversarial loss [55] and train the feature extractor  
193 to generate indistinguishable features that aim to fool the discriminators. This gives the following  
194 cross-domain adversarial loss:

$$\begin{aligned} \mathcal{L}_{CDAA} = & \sum_{i=1}^K \{ \mathbb{E}_{(x_{img}, x_{txt}) \sim (\mathcal{X}_{image}, \mathcal{X}_{text})} \log[D_i(G(x_{img}, x_{txt}))] \\ & + \mathbb{E}_{(x_{img}, x_{txt}) \sim (\mathcal{X}_{image}, \mathcal{X}_{text})} \log[1 - D_i(G(x_{img}, x_{txt}))] \}, \end{aligned} \quad (10)$$

195 where  $G$  denotes the feature extractor which includes the image and text encoder and the bi-linear  
196 pooling layer, we can see  $G(x_{img}, x_{txt})$  as the multi-modal feature  $\mathbf{f}$ , and  $D_i$  denotes the discrimina-  
197 tor belonging to source  $i$ .

## 198 2.5 M2CAN Learning

199 We can train a transferable sentiment classifier over the multi-modal feature space:  $f_t : \mathcal{F} \rightarrow \mathcal{Y}$ ,  
200 where  $\mathcal{F}$  is the space of multi-modal feature  $\mathbf{f}$ :

$$\mathcal{L}_{task} = -\mathbb{E}_{x_{img}, x_{txt}, y} \sim (\mathcal{X}_{image}^S, \mathcal{X}_{text}^S, Y_S) \left[ -\log P(y | f_t(G(x_{img}, x_{txt}))) \right]. \quad (11)$$

201 The final objective function of M2CAN is a weighted combination of different losses:

$$\mathcal{L}_{M2CAN} = \mathcal{L}_{task} + \lambda_1 \cdot \mathcal{L}_{CDAA} + \lambda_2 \cdot \mathcal{L}_{CMCA} + \lambda_3 \cdot \mathcal{L}_{CDCA}, \quad (12)$$

202 where  $\lambda_1, \lambda_2, \lambda_3$  are weights for different losses. This objective function can be optimized by solving  
203 the following min-max game:

$$f_t^* = \arg \min_{f_t} \min_G \max_{D_1, D_2} \mathcal{L}_{M2CAN}. \quad (13)$$

## 204 3 Experiments

205 Here we introduce the experimental settings and compare the sentiment classification results of  
206 M2CAN and several state-of-the-art DA approaches, followed by ablation study and visualization.

### 207 3.1 Experimental Settings

208 **Datasets.** Since we are the first to study the novel MS-MMDA setting and there is no specific dataset  
209 on this task, we evaluate our approach using a combined dataset, which consists of three public  
210 datasets on visual-textual sentiment: Yelp [12], Twitter [45], and MVSA [46]. We regard the three  
211 datasets as different domains since they follow different distributions. We create multiple MS-MMDA  
212 settings by taking each domain as *target* and the rest as *sources* in each setting.

213 The **Yelp domain** [12] contains customer-generated reviews of food services, *e.g.* restaurants,  
214 cafeterias, and dessert shops. In total, it has more than 44,000 reviews, including 244,000 images.  
215 Each review has a piece of textual comment, at least 3 images, and a score of sentiment polarity  
216 ranging from 1 to 5. We consider those reviews with scores of 1 and 2 as carrying negative sentiment,  
217 those with scores of 3 as carrying neutral sentiment, and those with scores of 4 and 5 as carrying  
218 positive sentiment. The **Twitter domain** [45] contains 50,000 user-generated tweets with images  
219 released on Twitter. Each tweet is composed of one textual review, several images, and a three-type  
220 sentiment label: negative, neutral, and positive. The **MVSA domain** [46] is also collected from  
221 Twitter. Similar to the Twitter domain, each tweet in MVSA consists of one textual review, several  
222 images, and a three-type sentiment label: negative, neutral, and positive. Specifically, each tweet is  
223 annotated by three experts. We abandon the tweets annotated with three different labels, and keep  
224 the tweets with at least two agreements. To balance the amount of samples in different domains, we  
225 randomly choose 15,000 samples as training set and 1,500 samples as test set for all domains.

226 **Evaluation Metrics.** Following [32, 29], we employ classification accuracy to evaluate the multi-  
227 modal sentiment classification results. Larger classification accuracy indicates better performance.

228 **Baselines.** We compare M2CAN with the following baselines: (1) **Source-only**, directly training on  
229 the source domains and testing on the target domain, which includes two settings: single-best, the  
230 best test accuracy on target among all source domains; source-combined, the target accuracy of the  
231 model trained on the combined source domain. (2) **Single-source MMDA methods**, including state-  
232 of-the-art approaches MMAN [23], MM-SADA [25], and xMUDA [26] trained with both single-best  
233 and source-combined settings. (3) **Multi-source MMDA methods**, including the state-of-the-art  
234 approach MDAN [32] and the proposed M2CAN. We also report the results of an oracle setting,  
235 where the model is both trained and tested on the target domain. We can view the oracle results as an  
236 upper bound for domain adaptation.

237 **Implementation Details.** For the image encoder, we use Resnet-50 [47]. For the text encoder, we  
238 use a 12-layer “bert-base-uncased” version BERT [48]. The weights for  $\mathcal{L}_{task}$ ,  $\mathcal{L}_{GAN}$ ,  $\mathcal{L}_{CMCA}$ , and  
239  $\mathcal{L}_{CDCA}$  are 1, 0.02, 0.02 and 0.05, respectively. We use a 2-layer multi-layer perceptron (MLP) to  
240 implement the lower-dimensional projection function  $g$ , and a fully-connected layer with activation  
241 function ReLU to implement both the discriminators and the task classifier. We use Adam [56] as the

Table 1: Comparison with the state-of-the-art DA methods on the combined dataset for visual-textual sentiment classification. All results are percentages. The best and second best classification accuracies trained on the source domains are emphasized with bold and underline respectively (same in Table 2).

Standards	Models	Yelp	Twitter	MVSA	Avg
Source-only	Source-combined (text only)	57.3	62.2	55.8	58.4
	Source-combined (text & image)	56.7	59.1	57.8	57.9
	Single-best (text & image)	56.9	61.8	57.2	58.6
Single-best MMDA	MMAN [23]	57.5	62.5	58.8	59.6
	MM-SADA [25]	58.1	<u>66.0</u>	58.3	60.8
	xMUDA [26]	<u>58.7</u>	64.3	57.6	60.2
Source-combined MMDA	MMAN [23]	55.8	64.2	60.8	60.3
	MM-SADA [25]	57.5	63.2	60.3	60.3
	xMUDA [26]	56.2	63.1	<b>62.0</b>	60.4
Multi-source MMDA	MDAN [32]	58.6	64.1	61.2	<u>61.3</u>
	<b>M2CAN (Ours)</b>	<b>62.5</b>	<b>67.9</b>	<u>61.6</u>	<b>64.0</b>
Oracle (train on target)		65.2	68.6	68.4	67.4

242 optimizer with a batch size of 8. The learning rate is 0.00002 for BERT and Resnet-50, and 0.0005  
 243 for the rest. All experiments are implemented in PyTorch and conducted on a machine with a Tesla  
 244 V100S-PCIE GPU with 32 GB memory. All implementation details are included in our source code.

### 245 3.2 Comparison With State-of-the-art

246 The performance comparisons between the proposed M2CAN and the baselines for visual-textual  
 247 sentiment classification, including source-only, single-source MMDA, and multi-source MMDA, are  
 248 shown in Table 1. From the results, we have the following observations:

249 (1) Without alleviating the domain shift between the source and target domains, both source-only  
 250 settings, *i.e.* single-best and source-combined, obtain poor classification accuracies, *i.e.* 58.6% and  
 251 57.8%, which are almost 10% worse than the oracle setting (67.4%). When setting Yelp and Twitter  
 252 as the target domain, it is clear that adding visual modality results in performance degradation for both  
 253 single-best and source-combined source-only settings as compared to using textual modality only,  
 254 *e.g.* 56.9% and 56.7 % vs. 57.3% on Yelp. This indicates that the large domain gap between source  
 255 and target domains results in severe interference between different modalities. These observations  
 256 motivate the research on domain adaptation.

257 (2) When directly applying to the MS-MMDA task, both single-best and source-combined MMDA  
 258 methods outperform the source-only setting. Since customers’ reviews vary a lot across domains,  
 259 features that are related to sentiment also differ. Therefore, these MMDA methods that can mitigate  
 260 the domain gap improve the sentiment classification results.

261 (3) Comparing the performances of source-combined and single-best MMDA methods, we can  
 262 find that naively performing single-source domain adaptation approaches on a combined dataset of  
 263 different sources might produce worse result (*i.e.* 60.3% of MM-SADA) than that on the best single  
 264 source (*i.e.* 60.8% of MM-SADA). This motivates our research on MS-MMDA.

265 (4) The proposed M2CAN performs the best (64.0%) among all adaptation settings. Compared to the  
 266 best results inside the source-only, single-best MMDA, source-combined MMDA, and multi-source  
 267 MMDA, M2CAN achieves 5.4%, 3.2%, 3.6%, and 2.7% performance gains, respectively. These re-  
 268 sults demonstrate that the proposed M2CAN model can achieve significantly better performance than  
 269 the state-of-the-art DA methods for visual-textual sentiment classification. The superior performance  
 270 of M2CAN benefits from the joint cross-modal alignment and cross-domain alignment.

### 271 3.3 Ablation Study

272 We conduct a series of ablation studies on the combined dataset to demonstrate the effectiveness of  
 273 different components of M2CAN. The results are shown in Table 2. First, we verify the necessity of  
 274 introducing extra modalities. Comparing the first two lines between with text only and with image

Table 2: Ablation study on different components of the proposed M2CAN on the combined dataset.

Models	Yelp	Twitter	MVSA	Avg
CDAA (text only)	57.2	62.8	58.1	59.4
CDAA (text & image)	58.5	64.2	<u>61.7</u>	61.5
CDAA + CMCA	61.1	64.5	61.6	62.4
CDAA + CDCA	<u>61.3</u>	<u>65.2</u>	<b>61.9</b>	<u>62.8</u>
<b>CDAA + CMCA + CDCA (M2CAN)</b>	<b>62.5</b>	<b>67.9</b>	61.6	<b>64.0</b>

275 and text for CDAA, we can see that after adding image, the average accuracy is improved by 2.1%,  
 276 demonstrating the effectiveness of introducing multiple modalities. Second, we investigate whether  
 277 it is necessary to construct the cross-domain adversarial alignment (CDAA). Comparing the first  
 278 two lines in Table 2 and Table 1, it is clear the performance is significantly improved (*e.g.* 61.5% vs.  
 279 57.9% when both image and text are used), demonstrating the effectiveness of adversarial alignment.  
 280 Third, we investigate the effectiveness of cross-modal contrastive alignment (CMCA). From the  
 281 third and second lines, we can see that compared to only using CDAA, adding CMCA achieves  
 282 0.9% performance gain on the average classification accuracy, which demonstrates the necessity  
 283 of the CMCA. Finally, we evaluate the influence of cross-domain contrastive alignment (CDCA).  
 284 Comparing CDAA vs. CDAA+CDCA and CDAA+CMCA vs. CDAA+CMCA+CDCA, we can  
 285 conclude that adding CDCA can further improve the performance, verifying that CDCA indeed  
 286 contributes to the adaptation task.

### 287 3.4 Visualization

288 In this section, we visualize the features of source and target samples before and after adaptation using  
 289 M2CAN. By using t-SNE [57] to reduce the dimensionality of samples, we plot the learned features  
 290 onto a 2-dimensional plane, with the results shown in Figure 2. Figure (a) represents the feature  
 291 representations before adaptation, while (b) represents the feature representations after adaptation by  
 292 M2CAN. Red represents source features and blue represents target features. As we can see, before  
 293 adaptation, source and target features can be obviously discriminated because of the existence of  
 294 domain gap; while after adaptation, we can hardly distinguish between source and target features.  
 295 Therefore, we can conclude that after adaptation the source and target features become more closely  
 296 aligned, which further demonstrates the effectiveness of M2CAN. Furthermore, we also plot the loss  
 297 curves in the training process. From Figure 3, we can observe that the different types of losses all  
 298 decline and converge through the training process.

### 299 3.5 Limitations

300 The proposed M2CAN works under the covariate shift and closed set assumptions [18, 22] with  
 301 labeled source data and unlabeled target data. When other domain shifts exist, such as label shift [22]  
 302 and category shift [39], we cannot guarantee satisfactory domain adaptation performances. As  
 303 stated in [22], there are many different domain adaptation settings, *i.e.* multiple target domains and  
 304 open-set labels. The proposed method does not explore such characteristics and thus cannot be  
 305 directly applied to these settings. Incorporating existing multi-source techniques, such as adding  
 306 discrimination between different sources [40], would improve the performance. Considering the  
 307 constraint of hardware resources, such as GPU memory, we did not exploit such techniques. Further,  
 308 because of the absence of datasets for multi-source multi-modal domain adaptation, we only verify  
 309 the effectiveness of M2CAN on a combined dataset with three different domains for visual-textual  
 310 sentiment classification. Extending the proposed method to other multi-modal domain adaptation  
 311 with multiple sources and exploring how to perform MS-MMDA when some sources contain few  
 312 labeled and sufficient unlabeled data remains our future work.

## 313 4 Conclusion

314 In this paper, we studied a novel and practical domain adaptation problem, *i.e.* multi-source multi-  
 315 modal domain adaptation (MS-MMDA), for visual-textual sentiment classification. The designed  
 316 multi-source multi-modal contrastive adversarial network (M2CAN) can learn domain-invariant multi-

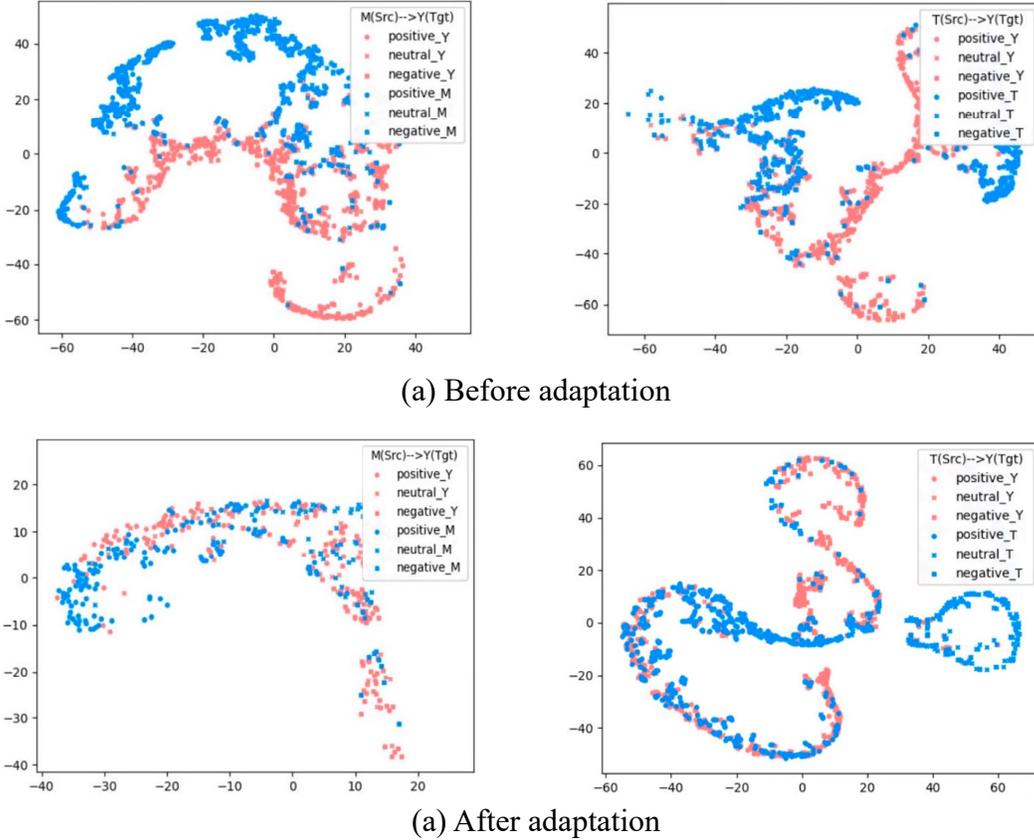


Figure 2: t-SNE visualization of multi-modal features before and after adaptation. Red represents source features and blue represents target features. We use Y, T, and M as abbreviations respectively for domains Yelp, Twitter, and MVSA for better visualization.

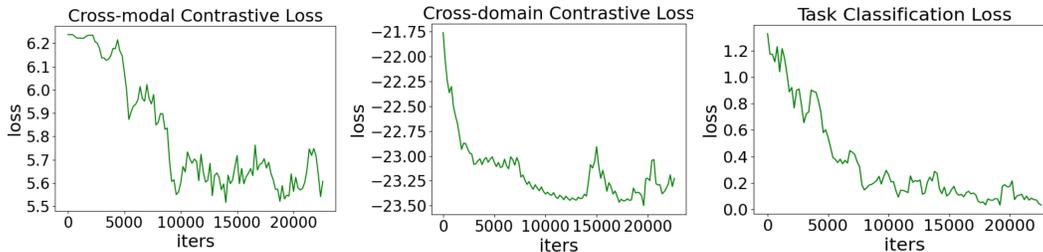


Figure 3: Visualization of different losses during the training process, including cross-modal contrastive loss, cross-domain contrastive loss and the task classification loss.

317 modal features by three different alignment strategies, *i.e.* cross-modal contrastive alignment within  
 318 each domain, cross-domain contrastive alignment for each modality, and cross-domain adversarial  
 319 alignment on the fused multi-modal representation. The cross-modal contrastive loss aligns visual and  
 320 textual features, pulling semantic-related sample pairs closer and pushing semantic-unrelated sample  
 321 pairs farther. The cross-domain contrastive loss together with domain adversarial loss bridge the  
 322 domain gap between the source and target domains while preserving the sentiment semantics through  
 323 contrastive learning and adversarial learning, respectively. Extensive experiments on a combined  
 324 dataset demonstrate the superiority of the proposed M2CAN as compared to the state-of-the-art DA  
 325 methods for visual-textual sentiment classification.

## 326 References

- 327 [1] Sicheng Zhao, Hongxun Yao, Yue Gao, Rongrong Ji, Wenlong Xie, Xiaolei Jiang, and Tat-Seng  
328 Chua. Predicting personalized emotion perceptions of social images. In *ACM International*  
329 *Conference on Multimedia*, pages 1385–1394, 2016.
- 330 [2] Jan Deriu, Aurelien Lucchi, Valeria De Luca, Aliaksei Severyn, Simon Müller, Mark Cieliebak,  
331 Thomas Hofmann, and Martin Jaggi. Leveraging large amounts of weakly supervised data for  
332 multi-language sentiment classification. In *International World Wide Web Conference*, pages  
333 1045–1052, 2017.
- 334 [3] Washington Luiz, Felipe Viegas, Rafael Alencar, Fernando Mourão, Thiago Salles, Dárlinton  
335 Carvalho, Marcos Andre Gonçalves, and Leonardo Rocha. A feature-oriented sentiment rating  
336 for mobile app reviews. In *International World Wide Web Conference*, pages 1909–1918, 2018.
- 337 [4] Lin Gong and Hongning Wang. When sentiment analysis meets social network: A holistic  
338 user behavior modeling in opinionated data. In *ACM SIGKDD International Conference on*  
339 *Knowledge Discovery and Data Mining*, pages 1455–1464, 2018.
- 340 [5] Niru Maheswaranathan, Alex H Williams, Matthew D Golub, Surya Ganguli, and David  
341 Sussillo. Reverse engineering recurrent networks for sentiment classification reveals line  
342 attractor dynamics. *Advances in Neural Information Processing Systems*, pages 15670–15679,  
343 2019.
- 344 [6] Sicheng Zhao, Yunsheng Ma, Yang Gu, Jufeng Yang, Tengfei Xing, Pengfei Xu, Runbo Hu, Hua  
345 Chai, and Kurt Keutzer. An end-to-end visual-audio attention network for emotion recognition  
346 in user-generated videos. In *AAAI Conference on Artificial Intelligence*, pages 303–311, 2020.
- 347 [7] Zhenpeng Chen, Sheng Shen, Ziniu Hu, Xuan Lu, Qiaozhu Mei, and Xuanzhe Liu. Emoji-  
348 powered representation learning for cross-lingual sentiment classification. In *International*  
349 *World Wide Web Conference*, pages 251–262, 2019.
- 350 [8] Quanzeng You, Liangliang Cao, Hailin Jin, and Jiebo Luo. Robust visual-textual sentiment  
351 analysis: When attention meets tree-structured recursive neural networks. In *ACM International*  
352 *Conference on Multimedia*, pages 1008–1017, 2016.
- 353 [9] Mohammad Soleymani, David Garcia, Brendan Jou, Björn Schuller, Shih-Fu Chang, and Maja  
354 Pantic. A survey of multimodal sentiment analysis. *Image and Vision Computing*, 65:3–14,  
355 2017.
- 356 [10] Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learn-  
357 ing: A survey and taxonomy. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,  
358 41(2):423–443, 2019.
- 359 [11] Jie Xu, Feiran Huang, Xiaoming Zhang, Senzhang Wang, Chaozhuo Li, Zhoujun Li, and  
360 Yueying He. Visual-textual sentiment classification with bi-directional multi-level attention  
361 networks. *Knowledge-Based Systems*, 178:61–73, 2019.
- 362 [12] Quoc-Tuan Truong and Hady W Lauw. Vistanet: Visual aspect attention network for multimodal  
363 sentiment analysis. In *AAAI Conference on Artificial Intelligence*, pages 305–312, 2019.
- 364 [13] Feiran Huang, Kaimin Wei, Jian Weng, and Zhoujun Li. Attention-based modality-gated  
365 networks for image-text sentiment analysis. *ACM Transactions on Multimedia Computing,*  
366 *Communications, and Applications*, 16(3):1–19, 2020.
- 367 [14] Eric Tzeng, Judy Hoffman, Trevor Darrell, and Kate Saenko. Simultaneous deep transfer across  
368 domains and tasks. In *IEEE International Conference on Computer Vision*, pages 4068–4076,  
369 2015.
- 370 [15] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei A  
371 Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *Internat-*  
372 *ional Conference on Machine Learning*, pages 1994–2003, 2018.

- 373 [16] Luyu Yang, Yogesh Balaji, Ser-Nam Lim, and Abhinav Shrivastava. Curriculum manager for  
374 source selection in multi-source domain adaptation. In *European Conference on Computer*  
375 *Vision*, pages 608–624, 2020.
- 376 [17] Antonio Torralba and Alexei A Efros. Unbiased look at dataset bias. In *IEEE Conference on*  
377 *Computer Vision and Pattern Recognition*, pages 1521–1528, 2011.
- 378 [18] Vishal M Patel, Raghuraman Gopalan, Ruonan Li, and Rama Chellappa. Visual domain  
379 adaptation: A survey of recent advances. *IEEE Signal Processing Magazine*, 32(3):53–69,  
380 2015.
- 381 [19] Shiliang Sun, Honglei Shi, and Yuanbin Wu. A survey of multi-source domain adaptation.  
382 *Information Fusion*, 24:84–92, 2015.
- 383 [20] Wouter Marco Kouw and Marco Loog. A review of domain adaptation without target labels.  
384 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.
- 385 [21] Sicheng Zhao, Bo Li, Colorado Reed, Pengfei Xu, and Kurt Keutzer. Multi-source domain  
386 adaptation in the deep learning era: A systematic survey. *arXiv preprint arXiv:2002.12169*,  
387 2020.
- 388 [22] Sicheng Zhao, Xiangyu Yue, Shanghang Zhang, Bo Li, Han Zhao, Bichen Wu, Ravi Krishna,  
389 Joseph E Gonzalez, Alberto L Sangiovanni-Vincentelli, Sanjit A Seshia, and Kurt Keutzer. A  
390 review of single-source deep unsupervised visual domain adaptation. *IEEE Transactions on*  
391 *Neural Networks and Learning Systems*, 2020.
- 392 [23] Xinhong Ma, Tianzhu Zhang, and Changsheng Xu. Deep multi-modality adversarial networks  
393 for unsupervised domain adaptation. *IEEE Transactions on Multimedia*, 21(9):2419–2431,  
394 2019.
- 395 [24] Fan Qi, Xiaoshan Yang, and Changsheng Xu. A unified framework for multimodal domain  
396 adaptation. In *ACM International Conference on Multimedia*, pages 429–437, 2018.
- 397 [25] Jonathan Munro and Dima Damen. Multi-modal domain adaptation for fine-grained action  
398 recognition. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
399 122–132, 2020.
- 400 [26] Maximilian Jaritz, Tuan-Hung Vu, Raoul de Charette, Emilie Wirbel, and Patrick Pérez. xmuda:  
401 Cross-modal unsupervised domain adaptation for 3d semantic segmentation. In *IEEE/CVF*  
402 *Conference on Computer Vision and Pattern Recognition*, pages 12605–12614, 2020.
- 403 [27] Huaiwen Zhang, Shengsheng Qian, Quan Fang, and Changsheng Xu. Multimodal disentangled  
404 domain adaption for social media event rumor detection. *IEEE Transactions on Multimedia*,  
405 2020.
- 406 [28] Massimiliano Mancini, Lorenzo Porzi, Samuel Rota Bulò, Barbara Caputo, and Elisa Ricci.  
407 Boosting domain adaptation by discovering latent domains. In *IEEE Conference on Computer*  
408 *Vision and Pattern Recognition*, pages 3771–3780, 2018.
- 409 [29] Jiang Guo, Darsh Shah, and Regina Barzilay. Multi-source domain adaptation with mixture  
410 of experts. In *Conference on Empirical Methods on Natural Language Processing*, pages  
411 4694–4703, 2018.
- 412 [30] Judy Hoffman, Mehryar Mohri, and Ningshan Zhang. Algorithms and theory for multiple-  
413 source adaptation. In *Advances in Neural Information Processing Systems*, pages 8246–8256,  
414 2018.
- 415 [31] Yitong Li, Michael Murias, Samantha Major, Geraldine Dawson, and David E Carlson. Extract-  
416 ing relationships by multi-domain matching. In *Advances in Neural Information Processing*  
417 *Systems*, pages 6799–6810, 2018.
- 418 [32] Han Zhao, Shanghang Zhang, Guanhang Wu, José MF Moura, Joao P Costeira, and Geoffrey J  
419 Gordon. Adversarial multiple source domain adaptation. In *Advances in Neural Information*  
420 *Processing Systems*, pages 8568–8579, 2018.

- 421 [33] Yongchun Zhu, Fuzhen Zhuang, and Deqing Wang. Aligning domain-specific distribution  
422 and classifier for cross-domain classification from multiple sources. In *AAAI Conference on*  
423 *Artificial Intelligence*, pages 5989–5996, 2019.
- 424 [34] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment  
425 matching for multi-source domain adaptation. In *IEEE International Conference on Computer*  
426 *Vision*, pages 1406–1415, 2019.
- 427 [35] Han Guo, Ramakanth Pasunuru, and Mohit Bansal. Multi-source domain adaptation for text  
428 classification via distancenet-bandits. In *AAAI Conference on Artificial Intelligence*, pages  
429 7830–7838, 2020.
- 430 [36] Sicheng Zhao, Guangzhi Wang, Shanghang Zhang, Yang Gu, Yaxian Li, Zhichao Song, Pengfei  
431 Xu, Runbo Hu, Hua Chai, and Kurt Keutzer. Multi-source distilling domain adaptation. In  
432 *AAAI Conference on Artificial Intelligence*, pages 12975–12983, 2020.
- 433 [37] Junfeng Wen, Russell Greiner, and Dale Schuurmans. Domain aggregation networks for  
434 multi-source domain adaptation. In *International Conference on Machine Learning*, pages  
435 10214–10224, 2020.
- 436 [38] Naveen Venkat, Jogendra Nath Kundu, Durgesh Kumar Singh, Ambareesh Revanur, and  
437 R Venkatesh Babu. Your classifier can secretly suffice multi-source domain adaptation. In  
438 *Advances in Neural Information Processing Systems*, 2020.
- 439 [39] Ziliang Chen, Pengxu Wei, Jingyu Zhuang, Guanbin Li, and Liang Lin. Deep cocktail networks.  
440 *International Journal of Computer Vision*, pages 1–24, 2021.
- 441 [40] Sicheng Zhao, Bo Li, Xiangyu Yue, Yang Gu, Pengfei Xu, Runbo Hu, Hua Chai, and Kurt  
442 Keutzer. Multi-source domain adaptation for semantic segmentation. In *Advances in Neural*  
443 *Information Processing Systems*, pages 7285–7298, 2019.
- 444 [41] Chuang Lin, Sicheng Zhao, Lei Meng, and Tat-Seng Chua. Multi-source domain adaptation for  
445 visual sentiment classification. In *AAAI Conference on Artificial Intelligence*, pages 2661–2668,  
446 2020.
- 447 [42] Sicheng Zhao, Yang Xiao, Jiang Guo, Xiangyu Yue, Jufeng Yang, Ravi Krishna, Pengfei Xu,  
448 and Kurt Keutzer. Curriculum cyclegan for textual sentiment domain adaptation with multiple  
449 sources. In *The Web Conference*, 2021.
- 450 [43] Jianzhong He, Xu Jia, Shuaijun Chen, and Jianzhuang Liu. Multi-source domain adaptation  
451 with collaborative learning for semantic segmentation. In *IEEE/CVF Conference on Computer*  
452 *Vision and Pattern Recognition*, 2021.
- 453 [44] Sicheng Zhao, Bo Li, Pengfei Xu, Xiangyu Yue, Guiguang Ding, and Kurt Keutzer. Madan:  
454 multi-source adversarial domain aggregation network for domain adaptation. *International*  
455 *Journal of Computer Vision*, pages 1–26, 2021.
- 456 [45] Qi Zhang, Jinlan Fu, Xiaoyu Liu, and Xuanjing Huang. Adaptive co-attention network for  
457 named entity recognition in tweets. In *AAAI Conference on Artificial Intelligence*, pages  
458 5674–5681, 2018.
- 459 [46] Teng Niu, Shiai Zhu, Lei Pang, and Abdulmotaleb El Saddik. Sentiment analysis on multi-view  
460 social data. In *International Conference on Multimedia Modeling*, pages 15–27, 2016.
- 461 [47] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
462 recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778,  
463 2016.
- 464 [48] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of  
465 deep bidirectional transformers for language understanding. In *Annual Conference of the North*  
466 *American Chapter of the Association for Computational Linguistics*, pages 4171–4186, 2019.
- 467 [49] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework  
468 for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020.

- 469 [50] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive  
470 predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- 471 [51] Guoliang Kang, Lu Jiang, Yi Yang, and Alexander G Hauptmann. Contrastive adaptation  
472 network for unsupervised domain adaptation. In *IEEE/CVF Conference on Computer Vision  
473 and Pattern Recognition*, pages 4893–4902, 2019.
- 474 [52] Joshua B Tenenbaum and William T Freeman. Separating style and content with bilinear models.  
475 *Neural Computation*, 12(6):1247–1283, 2000.
- 476 [53] Hamed Pirsiavash, Deva Ramanan, and Charless C Fowlkes. Bilinear classifiers for visual  
477 recognition. In *Advances in Neural Information Processing Systems*, pages 1482–1490, 2009.
- 478 [54] Jin-Hwa Kim, Kyoung-Woon On, Woosang Lim, Jeonghee Kim, Jung-Woo Ha, and Byoung-  
479 Tak Zhang. Hadamard product for low-rank bilinear pooling. In *International Conference on  
480 Learning Representations*, 2017.
- 481 [55] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
482 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural  
483 Information Processing Systems*, pages 2672–2680, 2014.
- 484 [56] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Interna-  
485 tional Conference on Learning Representations*, 2015.
- 486 [57] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine  
487 Learning Research*, 9(11):2579–2605, 2008.

## 488 Checklist

- 489 1. For all authors...
- 490 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
491 contributions and scope? [Yes]
- 492 (b) Did you describe the limitations of your work? [Yes] See Section 3.5.
- 493 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 494 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
495 them? [Yes]
- 496 2. If you are including theoretical results...
- 497 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 498 (b) Did you include complete proofs of all theoretical results? [N/A]
- 499 3. If you ran experiments...
- 500 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
501 perimental results (either in the supplemental material or as a URL)? [Yes] See the  
502 supplemental material.
- 503 (b) Did you specify all the training details (*e.g.*, data splits, hyperparameters, how they  
504 were chosen)? [Yes] See Section 3.1–Implementation Details.
- 505 (c) Did you report error bars (*e.g.*, with respect to the random seed after running experi-  
506 ments multiple times)? [N/A]
- 507 (d) Did you include the total amount of compute and the type of resources used (*e.g.*, type  
508 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 3.1–Implementation  
509 Details.
- 510 4. If you are using existing assets (*e.g.*, code, data, models) or curating/releasing new assets...
- 511 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 512 (b) Did you mention the license of the assets? [N/A]
- 513 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]  
514
- 515 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
516 using/curating? [N/A]

- 517 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
518 information or offensive content? [N/A]
- 519 5. If you used crowdsourcing or conducted research with human subjects...
- 520 (a) Did you include the full text of instructions given to participants and screenshots, if  
521 applicable? [N/A]
- 522 (b) Did you describe any potential participant risks, with links to Institutional Review  
523 Board (IRB) approvals, if applicable? [N/A]
- 524 (c) Did you include the estimated hourly wage paid to participants and the total amount  
525 spent on participant compensation? [N/A]