

CycleMatch: A Cycle-consistent Embedding Network for Image-Text Matching

Yu Liu^a, Yanming Guo^b, Li Liu^{b,c}, Erwin M. Bakker^a, Michael S. Lew^{a,*}

^a*Department of Computer Science, Leiden University, Leiden, 2333 CA, The Netherlands*

^b*College of System Engineering, National University of Defense Technology, Changsha, Hunan 410073, China*

^c*Center for Machine Vision and Signal Analysis, University of Oulu, Oulu, 8000, Finland.*

Abstract

In numerous multimedia and multi-modal tasks from image and video retrieval to zero-shot recognition to multimedia question and answering, bridging image and text representations plays an important and in some cases an indispensable role. To narrow the modality gap between vision and language, prior approaches attempt to discover their correlated semantics in a common feature space. However, these approaches omit the intra-modal semantic consistency when learning the inter-modal correlations. To address this problem, we propose cycle-consistent embeddings in a deep neural network for matching visual and textual representations. Our approach named as CycleMatch can maintain both inter-modal correlations and intra-modal consistency by cascading dual mappings and reconstructed mappings in a cyclic fashion. Moreover, in order to achieve a robust inference, we propose to employ two late-fusion approaches: average fusion and adaptive fusion. Both of them can effectively integrate the matching scores of different embedding features, without increasing the network complexity and training time. In the experiments on cross-modal retrieval, we demonstrate comprehensive results to verify the effectiveness of the proposed approach. Our approach achieves state-of-the-art performance on two well-known multi-modal datasets, Flickr30K and MSCOCO.

Keywords: Image-text matching, embedding, deep neural networks, late-Fusion inference

*Corresponding author

Email address: m.s.k.lew@liacs.leidenuniv.nl (Michael S. Lew)

1. Introduction

Nowadays, the explosive growth of multimedia data in social networks (*e.g.* image, video, text and audio) has triggered a massive amount of research activities in multi-modal understanding and reasoning. For instance, we can recognize a picture of a panda after hearing the description “black and white bears” without ever having seen one. This example demonstrates the cross-modal interaction between vision and language. These heterogeneous data offers us the opportunity to understand the world from diverse perspectives, while giving rise to the challenges of bridging different modalities. In this paper, we focus on the task of image-text matching, which aims to incorporate heterogeneous representations from visual and textual modalities. In practice, this task plays an essential role for a wide variety of tasks in the multimedia research, for examples, cross-modal retrieval [1, 2], visual question answering [3], zero-shot recognition [4] and video captioning [5].

The core issue with image-text matching is searching for an appropriate embedding space where related images and texts can be matched correctly. Driven by the great strides made by deep learning [6, 7, 8], recent research has been dedicated to exploring deep neural networks for learning powerful embedding features, in order to narrow the modality gap between visual and textual domains. These networks are typically composed of two branches for generating visual and textual embedding features in a common latent space, respectively [9, 10, 11, 12, 13]. Then, a similarity-based ranking loss is used to measure the latent embedding features. Latent embeddings can distill common semantic information about both the visual content and textual description. To directly match the similarities between vision and language, researchers further exploit dual embeddings by translating an input feature in the source space to be the feature in the target space [14, 15, 16, 17]. Both the latent and dual embeddings can capture inter-modal semantic correlations, however, they are limited in preserving intra-modal semantic consistency. Our motivation for this work is that: *A robust embedding method should be able to learn representations of both the source and target modalities.*

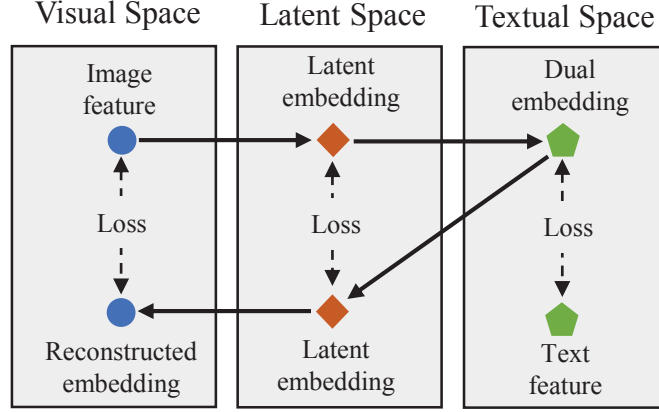
Inspired by the idea of cycle-consistent learning [18, 19], we propose cycle-consistent embeddings in an image-text matching network, which can incorporate both *inter-modal correlations* and *intra-modal consistency* for learning robust visual and textual embeddings. Figure 1 illustrates our embedding method by integrating three feature embeddings, including dual, reconstructed and latent embeddings. Specifically, it has two cycle branches, one starting from an image feature in the visual space and the other from a text feature in the textual space. For each branch, it first

accomplishes a dual mapping by translating an input feature in the source space to be a dual embedding in the target space. Inverse to the dual mapping, we then exploit a reconstructed mapping, with the aim of translating the dual embedding back to the source space. Moreover, we learn a latent space during the dual and reconstructed mappings and correlate the latent embeddings. In the three feature spaces, we compute their ranking losses to jointly optimize the whole embedding learning. Consequently, our visual-textual embedding method can learn not only *inter-modal mappings* (*i.e.* image-to-text and text-to-image), but also *intra-modal mappings* (*i.e.* image-to-image and text-to-text).

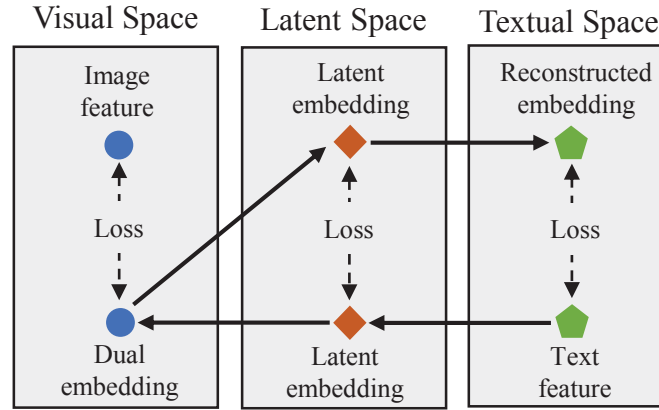
The contributions of this work are summarized as follows:

- We propose a novel deep cycle-consistent embedding network for image-text matching. Our approach called CycleMatch can cascade dual and reconstructed mappings together to maintain inter-modal correlations and intra-modal consistency. To our best knowledge, this is the first work to explore the usage of cycle consistency for solving the task of image-text matching.
- To improve the inference at the test stage, we present two late-fusion approaches to efficiently integrate the matching scores of multiple embedding features without increasing the training complexity.
- In the experiments, our cycle-consistency embedding outperforms traditional embeddings with considerable improvements for cross-modal retrieval on two multi-modal datasets, *i.e.* Flickr30K and MSCOCO. In addition, our results are competitive with the state-of-the-art performance on both datasets.

The rest of this paper is structured as follows. Related works are introduced in Section 2. Section 3 presents the details regarding the proposed CycleMatch. The late-fusion inference approaches are shown in Section 4. The experimental results are reported in Section 6. Finally, Section 7 summarizes the conclusions and discusses the future work.



(a) Image-to-text-to-image cycle



(b) Text-to-image-to-text cycle

Figure 1: Schematic pipeline of the proposed cycle-consistent embedding method. It is composed of two cycle branches: (a) image-to-text-to-image cycle and (b) text-to-image-to-text cycle. We first perform a dual mapping by transforming the input feature into the target feature space. Then, a reconstructed mapping is used to generate a reconstructed embedding in the source feature space. Moreover, we construct a latent space to correlate latent embeddings of the two mappings. The two branches share the mapping functions for transformations between three feature spaces, and can be trained jointly by optimizing the matching losses in the three feature spaces.

2. Related Work

Our work is related to image-text matching, deep visual-textual embedding and cycle-consistent learning.

2.1. Image-Text Matching

The problem of image-text matching has been studied by the multimedia community for decades. One typical solution is to unify heterogeneous representations into a latent embedding space, and then measure their similarity to ensure related pairs are more similar than unrelated ones. To be specific, Canonical Correlation Analysis (CCA) [20] is a classical and important embedding method, which can learn linear transformations to project two modalities into a latent space where their correlation is maximized. In addition, many variants [21, 22, 23, 24] are proposed to leverage the effectiveness of CCA. For example, kernel CCA [21] extended the classical linear CCA by learning non-linear transformations. Moreover, Gong *et al.* [25] integrated a third view with the two-view CCA using high-level image semantics, in order to gain a better separation for multi-modal data. Ranjan *et al.* [26] proposed a multi-label CCA approach by introducing multi-label information while learning the cross-modal subspaces. In practice, the integration of images and texts is a core issue for a variety of multi-modal applications [3, 4, 27, 28]. For example, Karaoglu *et al.* [29] proposed to detect words from images and then to combine the textual cues with the visual ones. Their method showed promising performance improvements for both place classification and logo retrieval. Similarly, Bai *et al.* [30] developed a unified and end-to-end trainable network, where the attention mechanism was further incorporated to better match the extracted textual and visual cues, to address the difficulties in fine-grained image classification.

2.2. Deep Visual-Textual Embedding

With the increasing progress of deep learning, research efforts have been made to CCA into deep neural networks [31, 32, 26, 33]. However, most deep CCA models rely on expensive decorrelation computations, which limit their generalization abilities at large-scale data. Alternatively, a number of recent approaches [34, 35, 12, 13, 36, 37] address the task by designing two-branch networks to embed visual and textual features into a common latent space, and then learn latent embeddings by optimizing a ranking loss between matched and unmatched image-text pairs. For instance, Wang *et al.* [9] built a simple and efficient matching network to preserve the structure relations

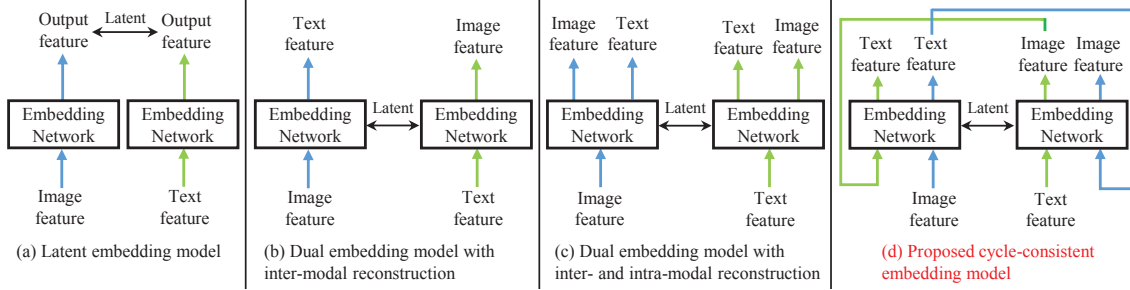


Figure 2: Conceptual illustration of variants of image-text matching models. (a) Latent embedding model. (b) Dual embedding model with inter-modal reconstruction. (c) Dual embedding model with inter-modal and intra-modal reconstruction. Note that each embedding network consists of two branches to output the image feature and text feature, separately. (d) Our cycle-consistent embedding model. The models in (b)(c)(d) also impose latent embeddings on hidden layers. Our model cascades the two embedding networks in a cyclic fashion, which can enhance interactions between two embedding networks.

between images and texts in the latent space. To associate image regions with words, the attention mechanism was integrated into visual-textual embedding models [10, 11]. In addition to the pairwise ranking loss, recent approaches [38, 39] leveraged extra loss functions to enhance the discrimination of the learned embedding features.

Another line of research [14, 15, 40, 41, 42] focuses on learning dual embeddings between two modalities, *e.g.* projecting visual features into the textual feature space and vice versa. Essentially, the dual embedding models are motivated by autoencoders. For instance, Feng *et al.* [14] proposed a correspondence cross-modal autoencoder model. 2WayNet [16] built the projections between two modalities and regularized them with Euclidean loss. Recently, the work of Gu *et al.* [17] utilized two generative models to synthesize grounded visual and textual representations. Also, Huang *et al.* [43] jointly modeled image-sentence matching and sentence generation. Note that, latent embeddings can be additionally used in the dual embedding models to enhance cross-modal relations.

In contrast to the above studies, our approach builds a reconstructed mapping upon the dual mapping, and generates cycle-consistent embeddings that are beneficial to the process of matching visual-textual representations. In Figure 2, we show the differences of our model from previous works.

2.3. Cycle-consistent Learning

There are a few papers exploring cycle consistency for diverse applications [18, 19, 44, 45, 46]. They are mainly motivated by the fact that, cycle-consistent learning is encouraged to produce additional feedback signals to improve the bi-directional translations. Specifically, He *et al.* [18] proposed a dual-learning mechanism based on deep reinforcement learning, where one agent was used to learn the primal task, *e.g.* English-to-French translation, and the other agent for the dual task, *e.g.* French-to-English translation. More recently, Zhu *et al.* [19] exploited cycle-consistent adversarial networks (CycleGAN), which combined a cycle-consistency loss with an adversarial loss [47] to perform unpaired image-to-image translations between two different visual domains. A similar idea was also presented in [48, 49]. Inspired by CycleGAN, several recent works have transferred the cycle-consistency loss to many supervised tasks [50, 51, 52].

Although prior works have shown the effectiveness of using cycle-consistent constraints for intra-modal domain mappings, yet in the context of cross-modal representation learning, its effectiveness has not been well investigated. In contrast to prior approaches that utilize cycle-consistent constraints within one modality (*e.g.* neural machine translation and image-to-image translation), our work is the first to extend the usage of cycle consistency for learning visual-textual embeddings. The work of Chen and Zitnick [53] is relevant to ours, as their model can both generate textual captions and reconstruct visual features given an image representation. However, their model lacks the inverse cycle mapping, *i.e.* text-to-image-to-text, which can be jointly learned in our model. Last but not least, these existing works did not consider matching latent embeddings during the cycle-consistent scheme.

3. Proposed Cycle-consistent Embeddings

In this section, we present the proposed CycleMatch model with cycle-consistent embeddings for matching visual and textual representations.

3.1. System Architecture

Figure 3 depicts an overview of the CycleMatch architecture. The entire network consists of three components: feature encoder, feature embedding and feature matching. First of all, given an input image I_i and text T_i , we employ individual feature encoders to extract the visual feature $\mathbf{v}_i = En_{img}(I_i)$ and textual feature $\mathbf{t}_i = En_{text}(T_i)$. Then, we develop several fully-connected (FC) layers

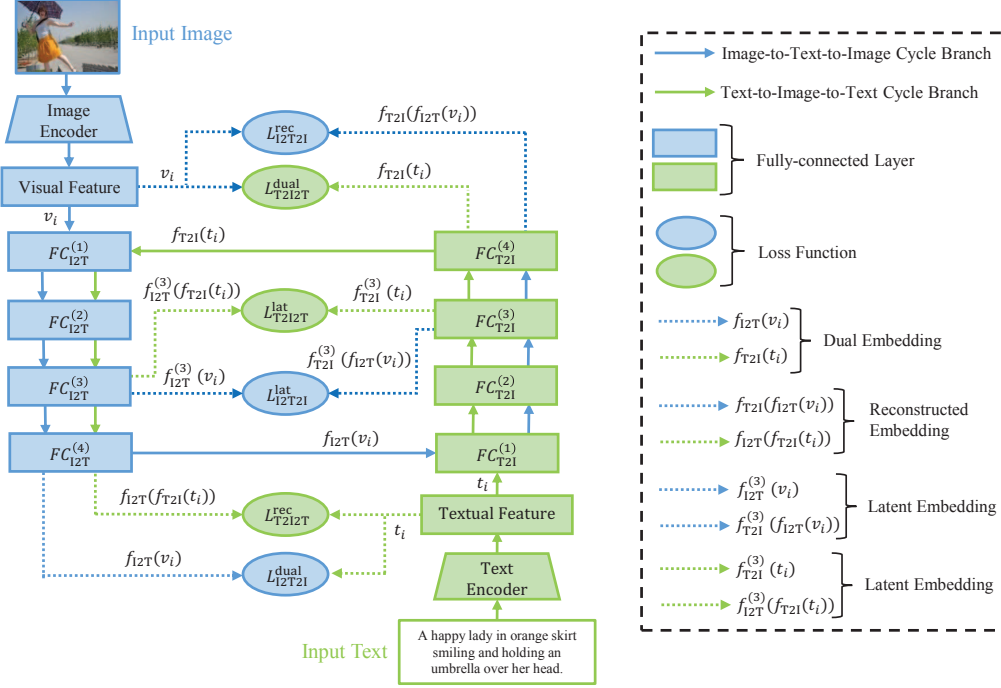


Figure 3: Overview of the proposed CycleMatch. It develops two cycle branches for visual-textual embeddings. For each branch, it is divided into two sub-branches from the *fourth* FC layer (*i.e.* $FC_{IT}^{(4)}$ and $FC_{TI}^{(4)}$). One sub-branch continues accomplishing the dual mapping to the target feature space, while the other sub-branch is used to perform the reconstructed mapping back to the source feature space. In this way, the cycle branches allow to jointly learn dual, reconstructed and latent embedding features. We can train the network end-to-end by optimizing several ranking loss functions simultaneously.

(*i.e.* $FC_{I2T}^{(j)}$) to perform the *Image-to-Text* (I2T) mapping and several other FC layers (*i.e.* $FC_{T2I}^{(j)}$) for the *Text-to-Image* (T2I) mapping. Let $f_{I2T}(\cdot)$ and $f_{T2I}(\cdot)$ represent the mapping functions for I2T and T2I, respectively. In addition, connecting FC_{I2T} and FC_{T2I} can form two cycle mappings between the visual and textual feature spaces. Specifically, given v_i , we first transform it to be $f_{I2T}(v_i)$ in the textual feature space and then learn its reconstructed feature $f_{T2I}(f_{I2T}(v_i))$ in the visual feature space. Moreover, we also correlate intermediate features derived from $FC_{I2T}^{(3)}$ and $FC_{T2I}^{(3)}$, so as to learn a latent feature space. Similarly, t_i is used to start another cycle mapping. In a nutshell, each cycle mapping can learn dual, reconstructed and latent embeddings in a cyclic fashion.

3.2. Formulation

Next, we will detail the above three embeddings and formulate their loss functions separately. The entire network contains two cycle-consistent embedding branches: one for *image-to-text-to-image* (I2T2I) mapping and the other for *text-to-image-to-text* (T2I2T) mapping. Here, we take the I2T2I mapping for an example.

3.2.1. Dual embedding

In a dataset collection with N image-text pairs, the input \mathbf{v}_i is fed into the first layer $FC_{\text{I2T}}^{(1)}$, where $i = 1, \dots, N$. By using the following layers $FC_{\text{I2T}}^{(j)}$ ($j = 2, 3, 4$), the network finally generates a dual embedding $f_{\text{I2T}}(\mathbf{v}_i)$ in the textual space, which has the same dimension as the ground-truth textual feature \mathbf{t}_i . Then, we normalize the two features and compute their similarity using the cosine distance

$$s(f_{\text{I2T}}(\mathbf{v}_i), \mathbf{t}_i) = \frac{f_{\text{I2T}}(\mathbf{v}_i) \cdot \mathbf{t}_i}{\|f_{\text{I2T}}(\mathbf{v}_i)\| \cdot \|\mathbf{t}_i\|}. \quad (1)$$

Notably, larger scores indicate more similar samples. During training, it is important to construct a number of negative pairs, in addition to the positive pair. Thereby, we search for the top K negative samples in a mini-batch for both $f_{\text{I2T}}(\mathbf{v}_i)$ and \mathbf{t}_i , which are denoted with $f_{\text{I2T}}(\mathbf{v}_{i,k}^-)$ and $\mathbf{t}_{i,k}^-$, respectively, where $k = 1, \dots, K$. To learn dual mappings, we need to employ a pairwise ranking loss function with respect to positive and negative pairs:

$$\begin{aligned} \mathcal{L}_{\text{I2T2I}}^{\text{dual}} = \sum_{i=1}^N \sum_{k=1}^K \Big\{ & \max \left[0, m - s(f_{\text{I2T}}(\mathbf{v}_i), \mathbf{t}_i) + s(f_{\text{I2T}}(\mathbf{v}_i), \mathbf{t}_{i,k}^-) \right] \\ & + \alpha \max \left[0, m - s(f_{\text{I2T}}(\mathbf{v}_i), \mathbf{t}_i) + s(f_{\text{I2T}}(\mathbf{v}_{i,k}^-), \mathbf{t}_i) \right] \Big\}, \end{aligned} \quad (2)$$

where m is a margin parameter that defines a threshold to constrain the positive and negative pairs. α adjusts the weights of the two loss terms. Ideally, the matched distance $s(f_{\text{I2T}}(\mathbf{v}_i), \mathbf{t}_i)$ should be smaller than any of the unmatched distances $s(f_{\text{I2T}}(\mathbf{v}_i), \mathbf{t}_{i,k}^-)$ and $s(f_{\text{I2T}}(\mathbf{v}_{i,k}^-), \mathbf{t}_i)$.

3.2.2. Reconstructed embedding

Despite the fact that the task in this work is about cross-modal matching, it is important as well to ensure intra-modal consistency, that is, related images (or texts) should have closer distances than unrelated ones. Hence, we explore reconstructed mappings to maintain the intra-modal semantic consistency, in addition to learning inter-modal correlations with dual mappings. We cascade

the dual and reconstructed mappings to form an intra-modal autoencoder and minimize the reconstruction error based on the ranking loss instead of the traditional Euclidean loss. Specifically, we feed $f_{I2T}(\mathbf{v}_i)$ into $FC_{T2I}^{(j)}$, to produce a reconstructed embedding feature $\tilde{\mathbf{v}}_i$ in the visual feature space with

$$\tilde{\mathbf{v}}_i = f_{T2I}(f_{I2T}(\mathbf{v}_i)) = f_{T2I} \circ f_{I2T}(\mathbf{v}_i). \quad (3)$$

The ranking loss for making the reconstructed embedding feature $\tilde{\mathbf{v}}_i$ match with the original visual feature \mathbf{v}_i can be written as follows

$$\begin{aligned} \mathcal{L}_{I2T2I}^{\text{rec}} = \sum_{i=1}^N \sum_{k=1}^K \Bigg\{ & \max \left[0, m - s(\tilde{\mathbf{v}}_i, \mathbf{v}_i) + s(\tilde{\mathbf{v}}_i, \mathbf{v}_{i,k}^-) \right] \\ & + \alpha \max \left[0, m - s(\tilde{\mathbf{v}}_i, \mathbf{v}_i) + s(\tilde{\mathbf{v}}_{i,k}^-, \mathbf{v}_i) \right] \Bigg\}. \end{aligned} \quad (4)$$

Since $\mathcal{L}_{I2T2I}^{\text{rec}}$ also has an effect on the parameters of $FC_{I2T}^{(j)}$, the reconstructed mappings can help to improve the learning of dual mappings as well.

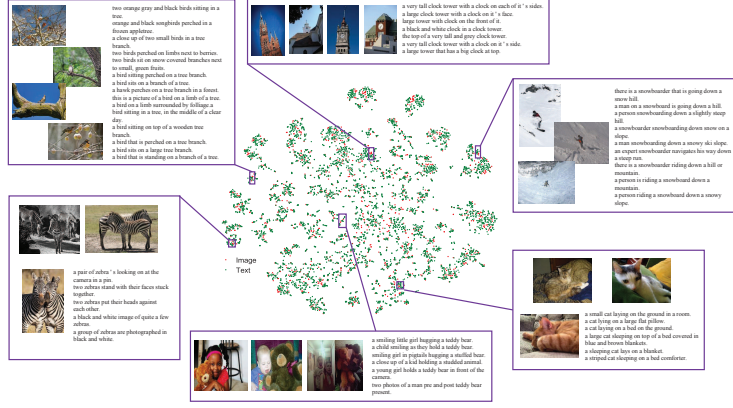
Moreover, we employ the t-SNE algorithm [54] to visualize our embedding features. Figure 4 shows the embedding maps with the test data from Flickr30K and MSCOCO, respectively. We show some original images and texts corresponding to the embedding features. First, the images and texts in each local window demonstrate high semantic correlations. In addition, these images themselves have similar visual content, and these texts themselves contain related descriptions. This observation is consistent with our motivation that a robust embedding method should be able to consider both inter-modal correlations and intra-modal consistency.

3.2.3. Latent embedding

Furthermore, we exploit a latent feature space to enhance the correlations between the dual and reconstructed mappings. Latent embeddings are able to distill common semantic information from visual and textual representations. Specifically, we make use of the intermediate representations from the third FC layers, *i.e.* $FC_{I2T}^{(3)}$ and $FC_{T2I}^{(3)}$. When \mathbf{v}_i passes through $FC_{I2T}^{(3)}$, we can extract an intermediate feature $f_{I2T}^{(3)}(\mathbf{v}_i)$. Also, the dual embedding $f_{I2T}(\mathbf{v}_i)$ passes through $FC_{T2I}^{(3)}$ to generate another intermediate feature $f_{T2I}^{(3)}(f_{I2T}(\mathbf{v}_i))$. The ranking loss for matching latent embeddings



(a) Flickr30K



(b) MSCOCO

Figure 4: Visualization of our embedding features. For each dataset, we pick 1000 images (red) and 5000 texts (green). Some images and texts corresponding to the embedding features are shown in local windows, from which we can observe not only correlations between cross-modal samples, but also relations between intra-modal samples.

thereby becomes

$$\begin{aligned}
\mathcal{L}_{I2T2I}^{\text{lat}} = & \sum_{i=1}^N \sum_{k=1}^K \left\{ \max \left[0, m - s(f_{I2T}^{(3)}(\mathbf{v}_i), f_{T2I}^{(3)}(f_{I2T}(\mathbf{v}_i))) \right. \right. \\
& + s(f_{I2T}^{(3)}(\mathbf{v}_i), f_{T2I}^{(3)}(f_{I2T}(\mathbf{v}_{i,k}^-))) \left. \right] \\
& + \alpha \max \left[0, m - s(f_{I2T}^{(3)}(\mathbf{v}_i), f_{T2I}^{(3)}(f_{I2T}(\mathbf{v}_i))) \right. \\
& \left. \left. + s(f_{I2T}^{(3)}(\mathbf{v}_{i,k}^-), f_{T2I}^{(3)}(f_{I2T}(\mathbf{v}_i))) \right] \right\}.
\end{aligned} \tag{5}$$

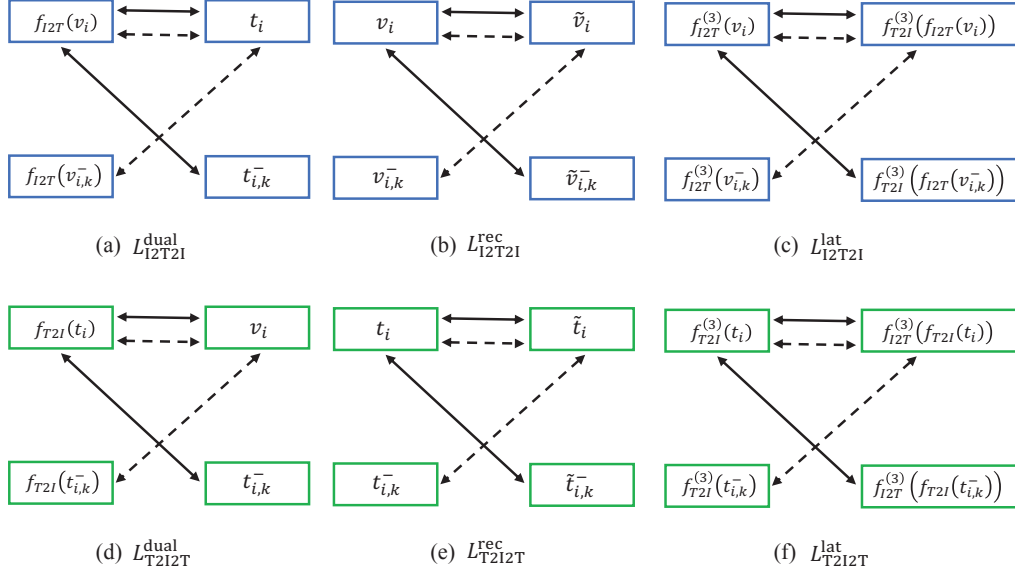


Figure 5: Conceptual illustration of loss functions for training CycleMatch. The first row includes the loss functions in the I2T2I cycle and the second row is for the T2I2T cycle.

Similar to the above I2T2I branch, it is straightforward to express the matching losses in the T2I2T branch, including $\mathcal{L}_{T2I2T}^{\text{dual}}$, $\mathcal{L}_{T2I2T}^{\text{rec}}$ and $\mathcal{L}_{T2I2T}^{\text{lat}}$. In Figure 5, we show the six loss functions for learning cycle-consistent embeddings.

3.2.4. Full objective

During training, we need to incorporate all the loss functions jointly. The full objective is to minimize the total loss:

$$\arg \min_{W_{I2T}, W_{T2I}} \mathcal{L}_{\text{total}} = \mathcal{L}_{I2T2I}^{\text{dual}} + \mathcal{L}_{I2T2I}^{\text{rec}} + \mathcal{L}_{I2T2I}^{\text{lat}} + \mathcal{L}_{T2I2T}^{\text{dual}} + \mathcal{L}_{T2I2T}^{\text{rec}} + \mathcal{L}_{T2I2T}^{\text{lat}}, \quad (6)$$

where W_{I2T} and W_{T2I} indicate the parameters in $FC_{I2T}^{(j)}$ and $FC_{T2I}^{(j)}$, respectively. They are unshared due to the specialization of two different modalities.

To demonstrate the effectiveness of our CycleMatch, we utilize the t-SNE [54] algorithm to visualize the embedding features learned in the visual, textual and latent feature spaces, separately. As shown in Figure 6, we randomly select 100 image-text pairs from the Flickr30K dataset [55]. From all the feature maps, we can visibly observe high similarities between two matched samples.

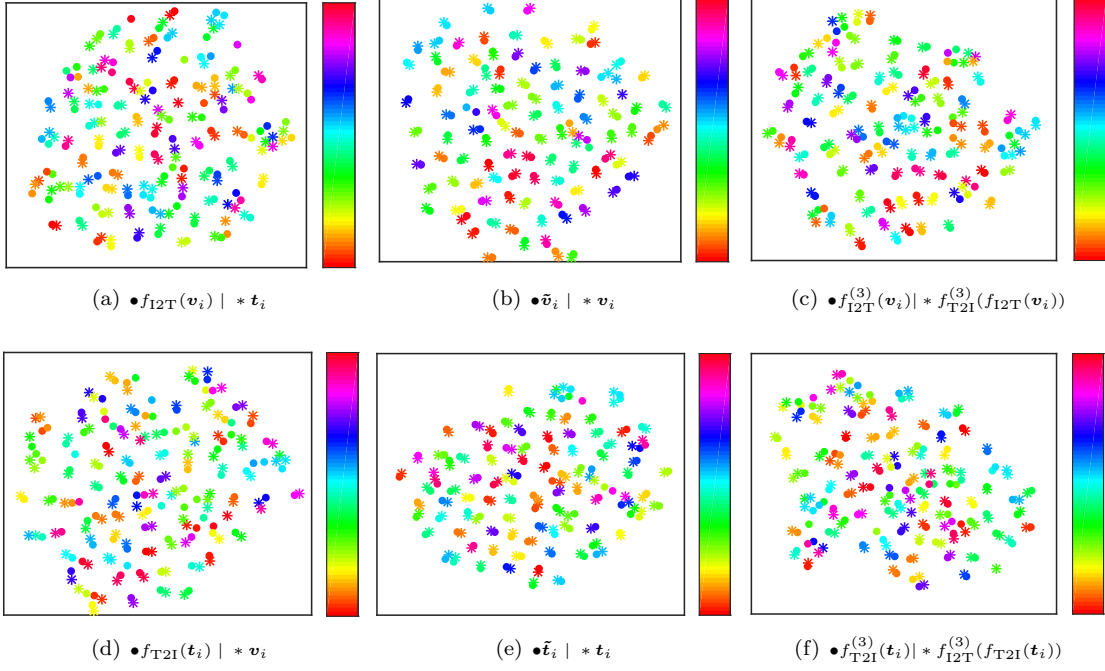


Figure 6: Visualization of our embedding features by using 100 image-text pairs in Flickr30K [55]. The first and second rows represent the embedding features learned in the I2T2I and T2I2T branches respectively. In each feature map, matched samples are shown with the same color. In (a)(d), the dual embedding features (\bullet) can match with the corresponding target features ($*$); In (b)(e), the reconstructed embedding features (\bullet) look closely similar to the source features ($*$). In (c)(f), the two latent embedding features (\bullet and $*$) can learn to correlate with each other as well.

4. Late-fusion Inference

By performing cycle-consistent embeddings, we can represent one sample with a set of three different features, for instance, $\{\mathbf{v}_i, f_{I2T}(\mathbf{v}_i), f_{I2T}^{(3)}(\mathbf{v}_i)\}$ for an image. Since the reconstructed embedding $\tilde{\mathbf{v}}_i$ and the other latent embedding $f_{T2I}^{(3)}(f_{I2T}(\mathbf{v}_i))$ are related to \mathbf{v}_i and $f_{I2T}^{(3)}(\mathbf{v}_i)$, we do not consider them for simplicity. Each of the three features can be used to measure an image-text matching score. Instead of using only one score, it is encouraged to leverage different scores together to achieve a more robust inference. This is driven by the late-fusion technique [56] in multimedia retrieval, which is a simple and efficient approach to combine the prediction scores of individual features. In this work, we present two effective late-fusion approaches, namely average fusion and adaptive fusion.

4.1. Average Fusion

Given a query image I_q , we extract three features $\{\mathbf{v}_q, f_{I2T}(\mathbf{v}_q), f_{I2T}^{(3)}(\mathbf{v}_q)\}$. Similarly, an arbitrary text T_i in the dataset can be described with $\{\mathbf{t}_i, f_{T2I}(\mathbf{t}_i), f_{T2I}^{(3)}(\mathbf{t}_i)\}$. We can compute three similarity scores between I_q and T_i :

$$\begin{cases} \text{visual score : } s^{(1)}(\mathbf{v}_q, \mathbf{t}_i) = s(\mathbf{v}_q, f_{T2I}(\mathbf{t}_i)), \\ \text{textual score : } s^{(2)}(\mathbf{v}_q, \mathbf{t}_i) = s(f_{I2T}(\mathbf{v}_q), \mathbf{t}_i), \\ \text{latent score : } s^{(3)}(\mathbf{v}_q, \mathbf{t}_i) = s(f_{I2T}^{(3)}(\mathbf{v}_q), f_{T2I}^{(3)}(\mathbf{t}_i)). \end{cases} \quad (7)$$

Then we combine the three scores to obtain an average fusion score as follows

$$s^{avg}(\mathbf{v}_q, \mathbf{t}_i) = \frac{\sum_{j=1}^3 s^{(j)}(\mathbf{v}_q, \mathbf{t}_i)}{3}. \quad (8)$$

It is similar to compute the fusion score $s^{avg}(\mathbf{t}_q, \mathbf{v}_i)$ in terms of a query text T_q .

4.2. Adaptive Fusion

To study the importance of different features, we further learn adaptive weights when combining the three scores. As suggested in [57], the score curve by using a superior feature can be sorted in an “L” shape, while the curve by using an inferior feature tends to gradually descend. In addition, the area under the curve can be used as an indicator to measure the weight of the corresponding feature. Driven by this observation, we can use the sorted score curves of the above three features to decide their weights. Specifically, we utilize each of the three features to compute the score curve of a query image I_q to all the text samples. Then, we sort the score curves and compute their areas with respect to the horizontal axis. In Figure 7, we show three sorted score curves for either a query image or text.

Our adaptive fusion method is inspired by the late fusion in [57], due to its parameter-free property and efficient computation. However, our method has two major differences from [57]. First, Zheng *et al.* [57] attempt to integrate different features, including BoW, Color and GIST features. In contrast, we construct a unified network to extract multiple embedding features, which have close relations to each other. Second, In [57], they use the total curve to compute the area.

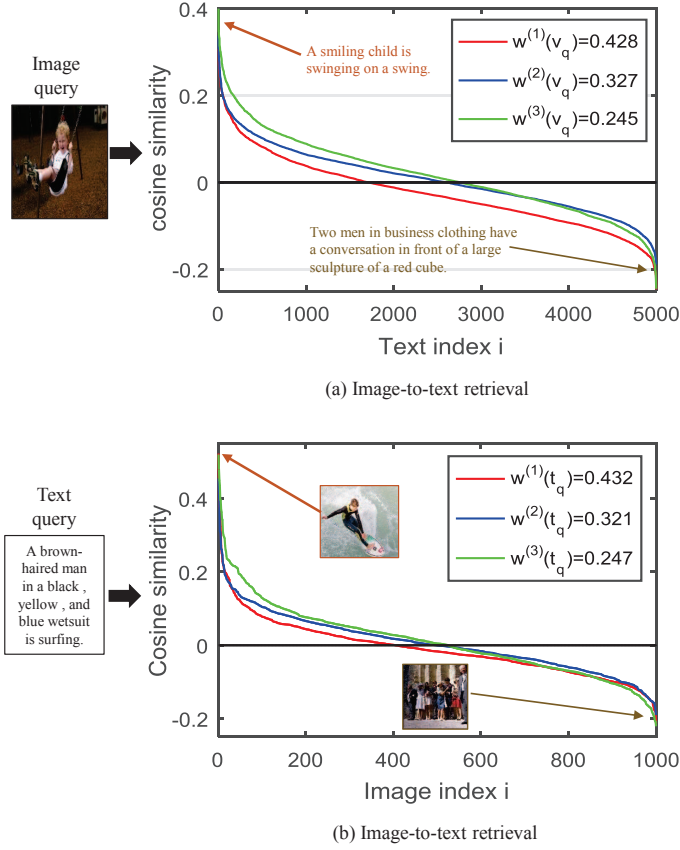


Figure 7: Illustration of the sorted score curves based on three different features. (a) For the query image, the first curve (in red) forms the smallest area above the X axis, so the corresponding feature (*i.e.* visual embedding feature) can have the largest weight (0.428). We show a matched text at the beginning of the curves and an unmatched text at the end of the curves. (b) Similarly, we demonstrate a text query example.

However, we compute only the positive area above the axis and omit the negative one ¹. This way can help to decrease the effect of long tails of the curves. For example in Figure 7 (b), the three curves have almost similar negative areas, based on which it is hard to distinguish the weights of the three features. Hence, adding the negative area with the positive one will narrow the gap of significance of different features and fail to learn robust adaptive weights. In the experiments, we show the advantage of our method over [57].

¹The similarity scores in this work are based on the cosine distance, ranging from -1 to 1.

Formally, the positive area associated with the j -th feature can be approximated by

$$area_+^{(j)}(\mathbf{v}_q) = \sum_{i=1}^N \max \left[0, s^{(j)}(\mathbf{v}_q, \mathbf{t}_i) \right]. \quad (9)$$

Smaller positive area means that the corresponding feature should have greater weights. Hence, the adaptive weights of I_q *w.r.t.* the three features can be expressed with

$$w^{(j)}(\mathbf{v}_q) = \frac{1}{area_+^{(j)}(\mathbf{v}_q)}. \quad (10)$$

In addition, we normalize the three weights to make sure $\sum_{j=1}^3 w^{(j)}(\mathbf{v}_q) = 1$. Finally, the adaptive fusion score for matching I_q and T_i becomes

$$s^{adt}(\mathbf{v}_q, \mathbf{t}_i) = \sum_j w^{(j)}(\mathbf{v}_q) \cdot s^{(j)}(\mathbf{v}_q, \mathbf{t}_i). \quad (11)$$

Likewise, we demonstrate a text query T_q in the right of Figure 7, and show its adaptive weights, $w^{(j)}(\mathbf{t}_q)$. Notice that our adaptive fusion approach can achieve specific weights for different query samples. It is an unsupervised and efficient manner without adding extra parameters and manual tuning. In the experiments, we analyze the effects of these two late-fusion approaches on the inference of cross-modal retrieval.

5. Discussion

Although the cycle-consistent idea has been adopted in many problems, it should not decrease the novelty of our work. In this section, we mainly aim to state our similarities and differences compared to the prior works like CycleGAN.

Similarities: Essentially, cycle-consistent learning is a variant of the auto-encoder model, which mainly aims to construct a cyclic mapping to reconstruct the input data. Both CycleGAN and CycleMatch are motivated by the idea of cycle-consistent learning, even though they focus on addressing different tasks.

Differences: Our proposed CycleMatch uses the idea of cycle-consistent learning, but it still has task-specific novelties and differences from CycleGAN.

- First, CycleGAN integrates a cycle-consistency loss with an adversarial loss to perform intra-modal representation learning, *i.e.* image-to-image translation between two image sets. In contrast, our CycleMatch is proposed to address the problem of cross-modal representation learning between image and text sets. In prior works, the effectiveness of cycle-consistent learning has not been well investigated in the context of cross-modal tasks. Our work is the first to extend cycle-consistent learning to address the task of image-text matching.
- Second, our reconstructed embedding is learned with the ranking loss, instead of the traditional Euclidean loss in CycleGAN. Notably, the ranking loss aims to reconstruct the relations among data samples rather than the original features. We find that the ranking loss is more suited for the matching task compared to the Euclidean loss.
- Third, CycleMatch is a novel network architecture that is different from CycleGAN. Notably, CycleMatch is not based on the GAN model. In addition, we consider the latent embedding representations, which are not taken into account in CycleGAN.
- Lastly, we contribute to proposing late-fusion inference in order to integrate multiple embedding features learned in the model. This robust and efficient inference is performed in the test stage and will not complicate the training procedure. The results in our experiments verify the effectiveness of the late-fusion inference. However, CycleGAN does not provide a robust inference in its test stage.

In summary, more and more papers [50, 51, 52, 58, 59] are making use of cycle-consistent learning to solve a variety of problems, such as domain adaptation, video retargeting and zero-shot learning. These works make promising contributions to the field, even though they are primarily or partially inspired by CycleGAN. In the future, we believe more related works will be encouraged in the field.

6. Experiments

First, we compare CycleMatch with various baseline models to verify its effectiveness. In addition, we present in-depth analysis on the two late-fusion approaches. Moreover, our results can be competitive with the state-of-the-art performance for cross-modal retrieval on two well-known datasets. Finally, we present additional ablation study on the effect of feature encoders and variance of test splits.

6.1. Experimental Setup

We introduce the dataset protocols, evaluation metrics, network details, training details and time complexity, involved in our experimental setup.

6.1.1. Dataset protocols

The experiments are performed on two well-known datasets. (1) Flickr30K [55] consists of 31,783 images and each image is associated with five different sentences. We use the dataset split of [60], namely 29,783 training images, 1,000 validation images and 1,000 test images. (2) MSCOCO [61] is one of the largest multi-modal datasets, which includes 82,783 training images and 40,504 validation images. We pick five ground-truth sentences for each image. 1,000 test images are selected from the validation set [60]. Notice that some works [9, 38, 17] merge the remaining validation images into the training set, to further increase the performance. However, we keep only using the original training set for fairness.

6.1.2. Evaluation metrics

For evaluating the performance of cross-modal retrieval, we adopt the common metric R@K, which measures the recall rate of a correctly retrieved ground-truth at top K retrieved candidates. Generally, K is set to 1, 5 and 10 for both image-to-text and text-to-image retrieval.

6.1.3. Network details

In terms of the image encoder, we employed the powerful ResNet-152 [8] pre-trained on the ImageNet dataset [62]. Besides, we recast the CNN model to its fully convolutional network (FCN) counterpart, which can capture rich region representations. The last layer of the FCN model is spatially averaged to generate a 2,048 dimensional visual representation. To extract the textual representation, we utilized the pre-trained RNN encoder proposed in [63]. It can represent one sentence with a 4,096 dimensional feature vector. Currently, we did not fine-tune the feature encoders during the training.

As for the two groups of four FC layers in CycleMatch (*i.e.* $FC_{12T}^{(j)}$ and $FC_{T2I}^{(j)}$), the channels of the first three layers are fixed as [2048, 512, 512]. Note that, $FC_{12T}^{(4)}$ should have the same dimension as the textual feature and $FC_{T2I}^{(4)}$ should be equal to the size of the visual feature.

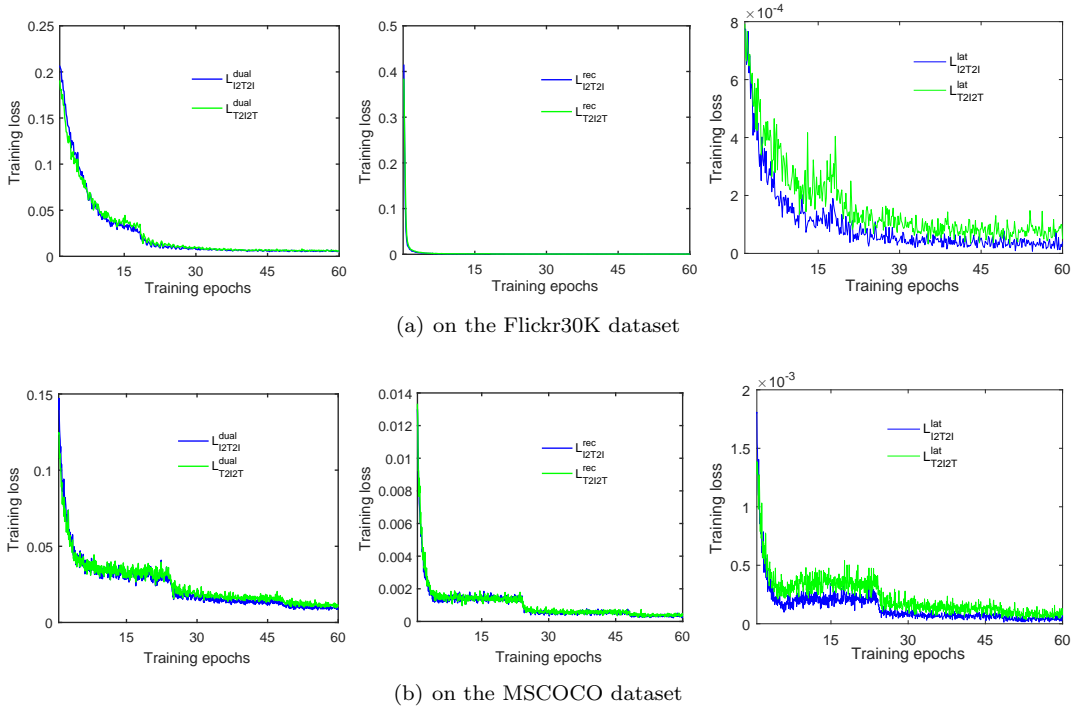


Figure 8: Illustration of training loss cost during training CycleMatch on the Flickr30K and MSCOCO datasets.

6.1.4. Training details

We implemented the proposed approach based on the Caffe library [64]. It is important to shuffle the training samples randomly during the data preparation stage. The hyper-parameters are evaluated on the validation set of each dataset. We trained the model using SGD with a mini-batch size of 500, a weight decay of 0.0005, a momentum of 0.9 and an initialized learning rate of 0.1. The learning rate is divided by 10 when the decrease in loss stabilizes. We set $\alpha = 2$ and $m = 0.1$ in all the experiments. The number of negative samples in each min-batch is 50. The whole training procedure terminates after 60 epochs for both datasets.

In Figure 8, we show the training loss of the six loss functions on the two datasets. It can be observed that the loss tend to converge during the training epochs.

6.1.5. Time complexity

We use the total loss in Eq. (6) to perform the training procedure. Each loss term is a simple and efficient ranking loss that is widely used in retrieval tasks. We used a Titan X card with 12

Table 1: Summary of various embedding methods for image-text matching.

Embedding methods	Main description
LatentMatch	a latent embedding model by matching $f_{I2T}^{(3)}(v_i)$ and $f_{T2I}^{(3)}(t_i)$.
DualMatch	a dual embedding model by learning two dual mappings: $I \rightarrow T$ and $T \rightarrow I$.
CycleMatch(w/o latent)	an ablation model without latent embeddings between dual and reconstructed mappings.
CycleMatch(I2T2I)	an ablation model with an I2T2I cycle branch and an $I \rightarrow T$ dual mapping.
CycleMatch(T2I2T)	an ablation model with a T2I2T cycle branch and a $T \rightarrow I$ dual mapping.
CycleMatch	a fully implemented model with two cycle branches.

Table 2: Comparison of different embedding approaches for cross-modal retrieval on Flickr30k and MSCOCO. Higher R@K numbers are better, where $K = 1, 5, 10$. The full CycleMatch method outperforms others on both datasets.

Method	Flickr30K dataset						MSCOCO dataset					
	Image to Text			Text to Image			Image to Text			Text to Image		
	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
LatentMatch	49.7	77.4	85.0	37.8	69.8	80.6	53.9	82.9	90.8	43.0	75.8	85.9
DualMatch	53.4	80.5	87.1	40.1	70.9	81.0	56.3	83.5	91.5	45.5	76.7	87.5
CycleMatch(w/o latent)	56.8	81.7	90.3	41.1	72.5	81.3	58.5	84.0	92.4	46.9	78.3	88.7
CycleMatch(I2T2I)	57.0	82.4	91.0	42.4	73.6	82.0	61.1	85.5	93.1	46.3	79.3	89.0
CycleMatch(T2I2T)	56.4	81.9	90.6	43.2	74.3	82.6	59.7	84.7	92.6	47.6	79.7	89.6
CycleMatch	57.8	83.3	90.9	43.2	74.8	83.8	60.5	86.3	93.7	47.2	80.3	90.4

GB to train all models in the experiments. For the full CycleMatch model, training required about 19 hours on the Flickr30K dataset and 47 hours on the MSCOCO dataset, respectively.

6.2. Comparisons with Baseline Approaches



To demonstrate the superiority of our approach, we designed several baseline models (see Table 1) based on the same network settings and training hyper-parameters as CycleMatch. In terms of inference, LatentMatch is evaluated with only the latent score. However, all the other models have both visual and textual scores. For consistency we utilize the average fusion approach to accomplish their inference. Table 2 reports the cross-modal retrieval performance of these models on both Flickr30K and MSCOCO. Overall, CycleMatch surpasses LatentMatch and DualMatch with significant improvements, and achieves overall superior performance over other variants of CycleMatch. In the next, we can report the results from several aspects.

6.2.1. Impact of reconstructed embeddings



First, we explain the benefit of constructing the cycle-consistent embeddings in our model. Primarily, cycle-consistent learning used in our model can benefit the dual embeddings. As can be seen in Figure 3, the reconstructed embeddings are built on top of the dual embeddings, therefore the reconstruction loss can help the training of the dual embeddings. The main difference between

DualMatch and CycleMatch(w/o latent) is that the latter model introduces a reconstructed mapping upon the traditional dual mapping. As reported in Table 2, the performance gap shows between DualMatch and CycleMatch(w/o latent) verifies the benefit of adding reconstructed embeddings in a cyclic fashion.

In addition to the above quantitative evaluation, we show image-to-text retrieval results as well to qualitatively compare the two methods. As shown in Figure 9, CycleMatch (w/o latent) can retrieve more accurate text descriptions than DualMatch, given the same query image. According to both quantitative and qualitative comparisons, it shows the improvements achieved by adding the cycle-consistency in our model.

Query image	Retrieved texts by DualMatch	Retrieved texts by CycleMatch (w/o latent)
	<p>A group of people sitting on a deck.</p> <p>A group of people are sitting outside a cafe drinking coffee and juice.</p> <p>A group of people sitting on a deck.</p> <p>People sitting outside a house enjoying wine.</p>	<p>A group of people sitting on a deck.</p> <p>A group of people sit on a deck.</p> <p>A group of people are sitting outside a cafe drinking coffee and juice.</p> <p>A group of people gathered out on a deck.</p>
	<p>A pair of skiers , a man and a woman , are climbing up a snowy and tree-lined hill.</p> <p>A woman with a blue jacket is posing for a picture while skiing down a mountain.</p> <p>Two people are skiing in the snowy mountains.</p> <p>A man wearing a blue jacket and a backpack is skiing through a snow covered forest.</p>	<p>A pair of skiers , a man and a woman , are climbing up a snowy and tree-lined hill.</p> <p>Two people smiling with skis with snow and trees everywhere.</p> <p>A woman with a blue jacket is posing for a picture while skiing down a mountain.</p> <p>Two people are skiing in the snowy mountains.</p>

(a) Flickr30K

Query image	Retrieved texts by DualMatch	Retrieved texts by CycleMatch (w/o latent)
	<p>a lady with lots of giant crabs cooking them on the grill.</p> <p>a lot of food that are in some baskets.</p> <p>the open air vendor is selling various kinds of seafood.</p> <p>a serving pot next to several trays of food.</p>	<p>a lady with lots of giant crabs cooking them on the grill.</p> <p>the open air vendor is selling various kinds of seafood.</p> <p>a lot of food that are in some baskets.</p> <p>a woman water crabs in some white trays and other foods.</p>
	<p>a baby next to a stuffed bear of some sort.</p> <p>a baby is sitting next to a stuffed bear.</p> <p>a smiling little girl hugging a teddy bear.</p> <p>the baby is sitting with a teddy bear.</p>	<p>a baby next to a stuffed bear of some sort.</p> <p>a baby is sitting next to a stuffed bear.</p> <p>a baby sleeping with teddy bear as big as he is.</p> <p>the baby is sitting with a teddy bear.</p>

(b) MSCOCO

Figure 9: Image-to-Text retrieval results on the datasets, (a) Flickr30K and (b) MSCOCO. The ground-truth descriptions are in green. By comparison, CycleMatch (w/o latent) achieves more accurate results than DualMatch.

Table 3: Evaluation on the effect of using different fully-connected layers on the latent embedding. The two-score adaptive fusion is used here. By comparison, $fc(3)$ is the best one for learning the latent embedding on most measurements.

Method	Flickr30K dataset						MSCOCO dataset					
	Image to Text			Text to Image			Image to Text			Text to Image		
	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
latent embedding with $FC^{(1)}$	57.2	82.1	90.6	42.1	73.5	83.3	59.5	85.0	93.4	46.9	79.3	89.4
latent embedding with $FC^{(2)}$	58.2	83.3	91.9	43.3	74.9	84.3	60.7	86.4	94.0	47.5	80.5	90.4
latent embedding with $FC^{(3)}$	58.6	83.6	91.6	43.6	75.3	84.2	61.1	86.8	94.2	47.9	80.9	90.9

6.2.2. Impact of latent embeddings

By comparing the results of CycleMatch and CycleMatch(w/o latent), we find that integrating the latent embeddings into CycleMatch brings further improvements over all metrics. For example, R@5 shows about 2% gains for both $I \rightarrow T$ and $T \rightarrow I$. Although using only latent embeddings (*i.e.* LatentMatch) is inferior to other models, it is beneficial to adopt them to improve other embedding methods like CycleMatch.

Moreover, we conduct an experiment below to test the effect of using different fully-connected (FC) layers on the latent embeddings. Apart from using the layer $FC^{(3)}$, we also test the latent embedding based on $FC^{(1)}$ or $FC^{(2)}$. In Table 3, we show the results by using three different FC layers. It can be seen that the both $FC^{(2)}$ and $FC^{(3)}$ features show better results than $FC^{(1)}$. Although $FC^{(1)}$ feature has more dimensions, its representation power is less than $FC^{(2)}$ and $FC^{(3)}$. One main reason is that $FC^{(1)}$ is the first layer in the network, but $FC^{(2)}$ and $FC^{(3)}$ are closer to the high-level semantics. In addition, $FC^{(3)}$ has slight improvements over $FC^{(2)}$. Based on these results, we decide to construct the latent embedding with the $FC^{(3)}$ features.

6.2.3. Impact of cycle branches

Both CycleMatch(I2T2I) and CycleMatch(T2I2T) can outperform LatentMatch and DualMatch, even though only one cycle-consistent embedding branch is used. By comparing these two models, CycleMatch(I2T2I) performs better for $I \rightarrow T$ retrieval, while CycleMatch(T2I2T) yields better results for $T \rightarrow I$ retrieval. When we incorporate the two cycle branches jointly, namely CycleMatch, it achieves overall superior performance over any single cycle branch on both datasets.

In addition to the R@K performance, we further present the matching scores computed by using our embedding features. To be specific, we randomly select 100 image-text pairs from the test set, and compute the similarity between one image and text. As shown in Figure 10, matched image-text pairs (with the same index) have greater similarity scores than unmatched ones.

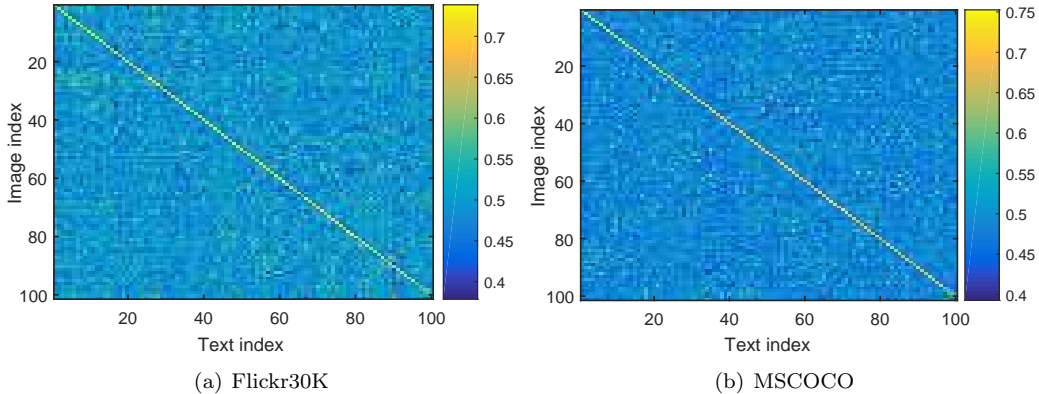


Figure 10: Similarity matrix of 100 image-text pairs from the test set. The related images and texts have the same index numbers. The diagonal line demonstrates high inter-modal correlations for matched image-text pairs. The original cosine scores are re-scaled to be $[0,1]$.

Table 4: Evaluation on the effect of different inference strategies on the R@K measurements. The two-score strategy based on the adaptive fusion achieves the best results (in bold face).

Inference method	Flickr30K dataset						MSCOCO dataset					
	Image to Text			Text to Image			Image to Text			Text to Image		
	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
One-score, without fusion	54.8	82.6	90.1	40.1	70.9	81.0	58.6	85.5	92.6	45.5	78.3	88.7
Two-score, average fusion	57.8	83.3	90.9	43.2	74.8	83.8	60.5	86.3	93.7	47.2	80.3	90.4
Two-score, adaptive fusion	58.6	83.6	91.6	43.6	75.3	84.2	61.1	86.8	94.2	47.9	80.9	90.9
Three-score, average fusion	57.4	83.5	91.0	43.2	74.7	83.9	59.7	86.0	94.0	46.9	80.6	89.8
Three-score, adaptive fusion	57.8	83.8	91.2	43.5	74.7	84.0	61.0	86.4	94.5	47.8	81.0	90.7

6.3. Analysis of Late-fusion Inference

Recall that CycleMatch contains visual, textual and latent scores for inference (Section 4). In this experiment, we compare three strategies to study the effect of two late-fusion inference approaches on the retrieval performance of CycleMatch. Specifically, the one-score strategy uses only a single visual score; the two-score strategy integrates visual and textual scores together; the three-score strategy combines all three scores by further adding the latent score. Table 4 reports the results of the three strategies. For the two-score and three-score strategies, we present the results of using the average and adaptive fusion, respectively. From the results, we can make the following observations:

- 1) The two-score strategy improves the one-score counterpart with 1%-3% gains. As the visual and textual scores match the samples in two different feature spaces, their complementary scores are able to improve the inference quality.

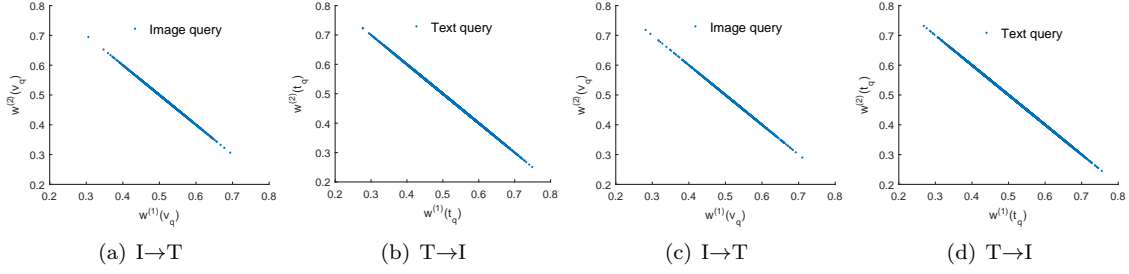


Figure 11: Visualization of adaptive weights for 1000 image queries and 5000 text queries on Flickr30K(a, b) and MSCOCO (c, d). Each dot in the maps is a query sample, having two weights for the adaptive fusion. Note that $w^{(1)}(\cdot) + w^{(2)}(\cdot) = 1$. The weights of query samples are mostly gathered between 0.4 and 0.6. It suggests that both visual and textual scores play an important role in the inference results.

2) The adaptive fusion outperforms the average one in terms of both two-score and three-score strategies. Although their performance gap over the R@K measurements is not significant, the adaptive fusion is an efficient method without imposing extra parameters and manual tuning. In addition, the inference time of the adaptive fusion is close to that of the average fusion.

3) The three-score strategy fails to achieve further improvements over the two-score one. We attribute this to the fact that, the latent score measures the similarity between $f_{I2T}^{(3)}(v_i)$ and $f_{T2I}^{(3)}(t_i)$. However, we do not use a direct matching loss between them during training CycleMatch. Although adding this latent score for inference will not bring further performance gains, learning the latent embeddings in CycleMatch is still important for improving the entire embedding procedure. As we discussed earlier, CycleMatch performs better than the variant without latent embeddings, namely CycleMatch(w/o latent).

As we can see, the two-score adaptive fusion achieves the best results. In Figure 11, we further present and analyze the two adaptive weights (*i.e.* $w^{(1)}(\cdot)$ and $w^{(2)}(\cdot)$), which are learned in the two-score adaptive fusion for visual and textual scores. Figure 11(a,b) and (c,d) shows the weights for Flickr30K and MSCOCO, respectively. For I2T retrieval, we illustrate the adaptive weights of 1000 image queries, namely $w^{(1)}(v_q)$ and $w^{(2)}(v_q)$; for T2I retrieval, we show all the weights of 5000 text queries, denoted as $w^{(1)}(t_q)$ and $w^{(2)}(t_q)$. Notice that, each dot in Figure 11 represents a query sample that learns individual weights based on its score curves. It can be seen that most samples have weights ranging from 0.4 to 0.6, which suggests that both visual and textual scores have an important impact on the inference results.

Comparison with the late fusion in [57]. This experiment is used to compare the results

Table 5: Comparison of two different methods for computing adaptive-fusion weights. The method by using only a positive area are better than that of using both positive and negative areas.

Inference method	Flickr30K dataset						MSCOCO dataset					
	Image to Text			Text to Image			Image to Text			Text to Image		
	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
Two-score, positive&negative areas	58.1	83.3	91.2	43.2	75.0	83.8	60.7	86.3	93.8	47.4	80.5	90.6
Two-score, positive area	58.6	83.6	91.6	43.6	75.3	84.2	61.1	86.8	94.2	47.9	80.9	90.9
Three-score, positive&negative areas	57.5	83.6	91.0	43.3	74.7	83.8	60.5	86.0	94.2	47.3	80.6	90.4
Three-score, positive area	57.8	83.8	91.2	43.5	74.7	84.0	61.0	86.4	94.5	47.8	81.0	90.7

Table 6: Comparison with the state-of-the-art approaches on Flickr30K. In addition, we present the image and text encoders used in these approaches. Our CycleMatch (the two-score adaptive fusion) achieves better results on R@K measurements (in boldface).

Method	Image encoder	Text encoder	Image to Text			Text to Image		
			R@1	R@5	R@10	R@1	R@5	R@10
DCCA [32]	AlexNet	TF-IDF	16.7	39.3	52.9	12.6	31.0	43.0
DVSA [34]	AlexNet	RNN	22.2	48.2	61.4	15.2	37.7	50.5
UVSE [36]	VGG-19	RNN	23.0	50.7	62.9	16.8	42.0	56.5
mCNN [35]	VGG-19	CNN	33.6	64.1	74.9	26.2	56.3	69.6
VQA-aware [65]	VGG-19	RNN	33.9	62.5	74.5	24.9	52.6	64.8
GMM-FV [33]	VGG-16	GMM+HGLMM	35.0	62.0	73.8	25.0	52.7	66.0
m-RNN [60]	VGG-16	RNN	35.4	63.8	73.7	22.8	50.7	63.1
RNN-FV [66]	VGG-19	RNN	35.6	62.5	74.2	27.4	55.9	70.0
HM-LSTM [13]	AlexNet	RNN	38.1	-	76.5	27.7	-	68.8
DSPE [9]	VGG-19	HGLMM	40.3	68.9	79.9	29.7	60.1	72.1
sm-LSTM [11]	VGG-19	RNN	42.5	71.9	81.5	30.2	60.4	72.3
VSE++ [67]	ResNet-152	RNN	43.7	-	82.1	32.2	-	72.1
DualCNN [38]	ResNet-152	ResNet-152	44.2	70.2	79.7	30.7	59.2	70.8
RRF-Net [12]	ResNet-152	HGLMM	47.6	77.4	87.1	35.4	68.3	79.9
2WayNet [16]	VGG-16	GMM+HGLMM	49.8	67.5	-	36.0	55.6	-
DAN [10]	ResNet-152	RNN	55.0	81.8	89.0	39.4	69.2	79.1
CycleMatch (Ours)	ResNet-152	RNN	58.6	83.6	91.6	43.6	75.3	84.2

of our adaptive fusion and the one in [57]. Recall that our method computes only the positive area above the axis, while the method in [57] considers both positive and negative areas. As reported in Table 5 below, the results with only a positive area are better in the context of both two-score and three-score fusion cases, even though the performance gap between the two methods is slight.

6.4. Comparisons with State-of-the-art Approaches

In Table 6 and Table 7, we present a comprehensive comparison with previous papers where they reported the cross-modal retrieval performance on Flickr30K and MSCOCO. It can be seen that our CycleMatch (the two-score adaptive fusion) outperforms recent state-of-the-art approaches [10, 12, 67] with promising improvements on both datasets. It is worth noting that these approaches

Table 7: Comparison with the state-of-the-art approaches on MSCOCO. In addition, we present the image and text encoders used in these approaches. Our CycleMatch (the two-score adaptive fusion) outperforms other approaches by achieving promising results (in boldface).

Method	Image encoder	Text encoder	Image to Text			Text to Image		
			R@1	R@5	R@10	R@1	R@5	R@10
STV [68]	VGG-19	RNN	33.8	67.7	82.1	25.9	60.0	74.6
DVSA [34]	AlexNet	RNN	38.4	69.9	80.5	27.4	60.2	74.8
GMM-FV [33]	VGG-16	GMM+HGLMM	39.4	67.9	80.9	25.1	59.8	76.6
m-RNN [60]	VGG-16	RNN	41.0	73.0	83.5	29.0	42.2	77.0
RNN-FV [66]	VGG-19	RNN	41.5	72.0	82.9	29.2	64.7	80.4
BiLSTM-Max [63]	ResNet-101	RNN	42.6	75.3	87.3	33.9	69.7	83.8
mCNN [35]	VGG-19	CNN	42.8	73.1	84.1	32.6	68.6	82.8
UVSE [36]	VGG-19	RNN	43.4	75.7	85.8	31.0	66.7	79.9
HM-LSTM [13]	AlexNet	RNN	43.9	-	87.8	36.1	-	86.7
order-embeddings [69]	VGG-19	RNN	46.7	-	88.9	37.9	-	85.9
DSPE [9]	VGG-19	HGLMM	50.1	79.7	89.2	39.6	75.2	86.9
VQA-aware [65]	VGG-19	RNN	50.5	80.1	89.7	37.0	70.9	82.9
DualCNN [38]	ResNet-50	ResNet-50	52.2	80.4	88.7	37.2	69.5	80.6
sm-LSTM [11]	VGG-19	RNN	53.2	83.1	91.5	40.7	75.8	87.4
2WayNet [16]	VGG-16	GMM+HGLMM	55.8	75.2	-	39.7	63.3	-
RRF-Net [12]	ResNet-152	HGLMM	56.4	85.3	91.5	43.9	78.1	88.6
VSE++ [67]	ResNet-152	RNN	58.3	-	93.3	43.6	-	87.8
CycleMatch (Ours)	ResNet-152	RNN	61.1	86.8	94.2	47.9	80.9	90.9

employ different feature encoders that have a significant influence on the performance. For a clear comparison, we further list the image and text encoders used in these approaches. In the following experiments, we will study the effect of different feature encoders on the performance of CycleMatch.

To boost the performance, recent several approaches [38, 67, 17, 43] further fine-tune the image encoders during training their models. Their results with fine-tuning the image encoders achieve better performance on MSCOCO than Flickr30K. We should know that it is feasible to fine-tune the image encoders while training our CycleMatch, which can help to further improve our results. In addition, the fine-tuning process will maintain the findings we mentioned as above. More importantly, our results on the Flickr30K dataset can even compete with the fine-tuned results in [38, 67, 17, 43]. On the MSCOCO dataset, the fine-tuned approaches [38, 67, 17, 43] merge the validation images into the training set to further increase the performance. However, we still use the original training set for a fair comparison with other prior approaches. Notice that, the accuracies of image-to-text retrieval are higher than those of text-to-image retrieval, because one image is annotated with several texts but one text matches with only one image.

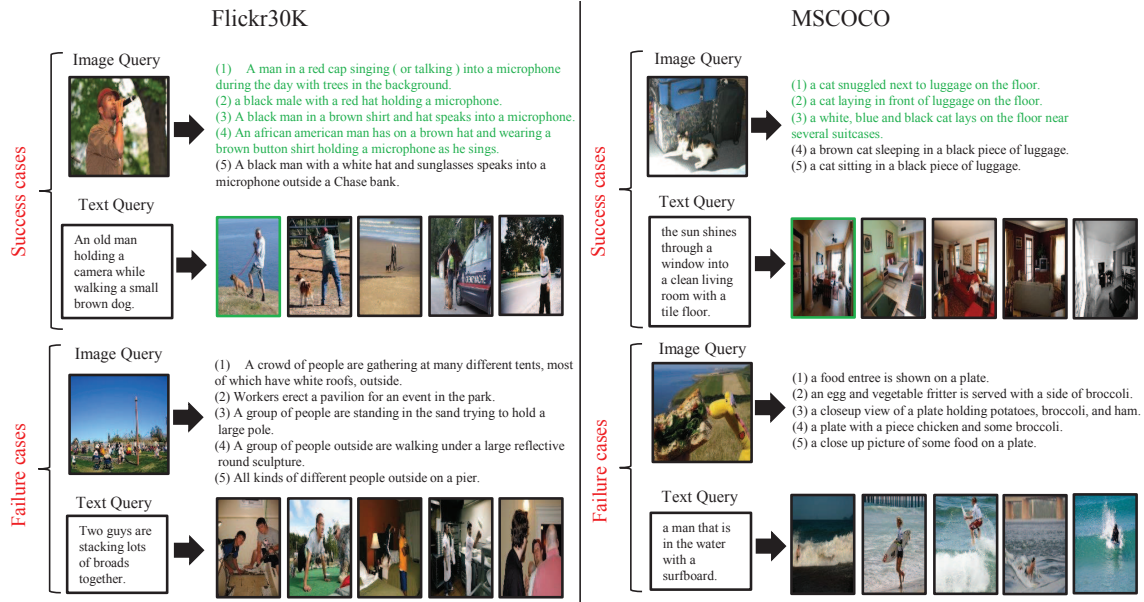


Figure 12: Qualitative results of our CycleMatch on Flickr30K and MSCOCO. Given one query, the top-5 candidates are retrieved. In the success cases, the correct matches are highlighted with green. In the failure cases, our method can still retrieve some reasonable candidates related to the query.

In addition to the quantitative evaluation, we present our image-to-text and text-to-image retrieval examples in Figure 12, which includes both success and failure cases. For each query sample, the top-5 candidates are retrieved, of which the ground-truth samples are highlighted in green. We notice that, the retrieved candidates are semantically related to the query sample in some extent, even for the failure cases.

6.5. Effect of Feature Encoders

As shown in Figure 3, we extract visual and textual features from off-the-shelf feature encoders. The proposed CycleMatch can be compatible with diverse feature encoders, but it is still encouraged to study the effect of different feature encoders on the performance. We report the results in Table 8.

Considering the image encoders, we use the VGG-19 and ResNet-152 models to extract the visual features and compare their results. We can see that, ResNet-152 has a considerable improvements over VGG-19 on all measurements, especially for R@1 accuracies. This shows the benefit of using more powerful CNN models for improving the visual embeddings. In addition, the feature dimension with ResNet-152 (*i.e.* 2,048) is lower than that with VGG-19 (*i.e.* 4,096). Therefore, in this work

Table 8: Evaluation on the effect of different feature encoders on the performance of CycleMatch. By comparison, ResNet-152 is a superior image encoder and RNN is a more powerful text encoder.

Image encoder	Text encoder	Flickr30K						MSCOCO					
		Image to Text			Text to Image			Image to Text			Text to Image		
		R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
Effect of image encoders													
VGG-19	RNN	51.4	80.6	88.1	38.5	71.0	81.3	55.1	83.5	91.3	43.7	76.7	88.4
ResNet-152	RNN	58.6	83.6	91.6	43.6	75.3	84.2	61.1	86.8	94.2	47.9	80.9	90.9
Effect of text encoders													
ResNet-152	word2vec	48.1	78.7	87.4	37.7	70.8	81.1	55.9	83.8	91.8	44.7	79.1	87.7
ResNet-152	HGLMM	54.5	81.6	90.9	41.3	73.1	82.8	58.4	85.5	93.4	46.2	80.3	89.4
ResNet-152	RNN	58.6	83.6	91.6	43.6	75.3	84.2	61.1	86.8	94.2	47.9	80.9	90.9

Table 9: Evaluation on the effect of fine-tuning the image encoder during training CycleMatch.

Image encoder	Text encoder	Flickr30K						MSCOCO					
		Image to Text			Text to Image			Image to Text			Text to Image		
		R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
VGG-19 without fine-tune		51.4	80.6	88.1	38.5	71.0	81.3	55.1	83.5	91.3	43.7	76.7	88.4
VGG-19 with fine-tune		54.8	83.1	90.5	42.3	74.8	84.5	60.2	87.3	94.0	49.3	81.2	91.8

we take the ResNet-152 model as the preferable image encoder.

In terms of the text encoders, we test another two encoders apart from the RNN encoder. The first one is word2vec [70], which describes each word in the sentence with a 300-dimensional feature vector. We then compute the average of all the word features to represent the sentence feature. The second one is an expensive representation based on the Hybrid Gaussian-Laplacian mixture model (HGLMM) [33]. Specifically, HGLMM computes a 18,000-dimension feature vector with 30 centers (*i.e.* 300*30*2). Similar to [9], we further reduce it to a 6,000-dimension feature vector in order to decrease the training complexity. As shown in Table 8, the RNN encoder is more powerful than both word2vec and HGLMM. In addition, the feature dimension based RNN (*i.e.* 4,096) is feasible and practical during training CycleMatch.

6.6. Effect of Fine-tuning Image Encoders

In this experiment, we perform the fine-tuning (ft) process for the VGG-19 image encoder. The results in Table 9 show considerable improvements for all R@K measurements. Similarly, fine-tuning ResNet-152 can bring further improvements as well, while it is out of the scope in our work.

Table 10: Evaluation on the effect of different test splits on the performance of CycleMatch. The results on the MSCOCO dataset show that CycleMatch can achieve high mean accuracy and low standard deviation.

Image encoder	Text encoder	Image to Text			Text to Image		
		R@1	R@5	R@10	R@1	R@5	R@10
ResNet-152	RNN	60.23 \pm 1.46	88.08 \pm 1.19	94.88 \pm 0.77	47.73 \pm 0.91	81.89 \pm 0.88	91.41 \pm 0.62

6.7. Variance of Test Splits

For a fair comparison, we employ the standard data split including 1,000 test images that are captured from the validation set [60, 12]. However, no prior work has studied the effect of using different test splits on the retrieval performance. To this end, we perform 100 times of evaluations on the MSCOCO dataset. For each evaluation, we randomly select 1,000 images from the validation set and test the results with the proposed CycleMatch. As shown in Table 10, our results show high mean accuracy and low standard deviation. This reveals the proper stability of our approach for cross-modal retrieval. It is worth mentioning that we cannot conduct this experiment on the Flickr30K dataset, as its test set (*i.e.* including only 1000 images) has been already fixed.

7. Conclusions

In this paper, we have developed a novel embedding method for the multi-modal task of matching visual and textual representations. We proposed cycle-consistent embeddings to learn both intra-modal correlations and intra-modal consistency. Our approach taking advantage of multiple embedding techniques is able to outperform any single embedding method. The experimental results have demonstrated the superiority of our method over other embedding methods. In addition, we have presented two simple and efficient late-fusion approaches to increase the inference quality. The late-fusion inference can integrate different matching scores together without increasing the training complexity. Finally, our approach has shown state-of-the-art performance for cross-modal retrieval on Flickr30K and MSCOCO. In the future, we will take into account local relations when matching images and sentences, for example, semantic correlations between visual regions and phases. One potential solution is to exploit the attention mechanism to localize the objects corresponding to the phase description.

Acknowledgments

This work was supported by the LIACS Media Lab at Leiden University under Grant 2006002026

and the National Natural Science Foundation of China under Grant 61872379. We are also grateful to the support of NVIDIA with the donation of GPU cards.

References

- [1] D. Rafailidis, S. Manolopoulou, P. Daras, A unified framework for multimodal retrieval, *Pattern Recognition* 46 (12) (2013) 3358–3370.
- [2] V. E. Liong, J. Lu, Y.-P. Tan, Cross-Modal Discrete Hashing, *Pattern Recognition* 79 (2018) 114–129.
- [3] M. Malinowski, M. Rohrbach, M. Fritz, Ask Your Neurons: A Neural-Based Approach to Answering Questions About Images, in: *IEEE ICCV*, 1–9, 2015.
- [4] S. Reed, Z. Akata, H. Lee, B. Schiele, Learning Deep Representations of Fine-Grained Visual Descriptions, in: *IEEE CVPR*, 49–58, 2016.
- [5] L. Gao, Z. Guo, H. Zhang, X. Xu, H. T. Shen, Video Captioning With Attention-Based LSTM and Semantic Consistency, *IEEE Transactions on Multimedia* 19 (9) (2017) 2045–2055.
- [6] A. Krizhevsky, I. Sutskever, G. E. Hinton, ImageNet Classification with Deep Convolutional Neural Networks, in: *NIPS*, 1097–1105, 2012.
- [7] K. Simonyan, A. Zisserman, Very Deep Convolutional Networks for Large-Scale Image Recognition, in: *ICLR*, 2015.
- [8] K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, in: *IEEE CVPR*, 770–778, 2016.
- [9] L. Wang, Y. Li, S. Lazebnik, Learning Deep Structure-Preserving Image-Text Embeddings, in: *IEEE CVPR*, 5005–5013, 2016.
- [10] H. Nam, J.-W. Ha, J. Kim, Dual Attention Networks for Multimodal Reasoning and Matching, in: *IEEE CVPR*, 299–307, 2017.
- [11] Y. Huang, W. Wang, L. Wang, Instance-Aware Image and Sentence Matching With Selective Multimodal LSTM, in: *IEEE CVPR*, 2310–2318, 2017.
- [12] Y. Liu, Y. Guo, E. M. Bakker, M. S. Lew, Learning a Recurrent Residual Fusion Network for Multimodal Matching, in: *IEEE ICCV*, 4107–4116, 2017.
- [13] Z. Niu, M. Zhou, L. Wang, X. Gao, G. Hua, Hierarchical Multimodal LSTM for Dense Visual-Semantic Embedding, in: *IEEE ICCV*, 1881–1889, 2017.
- [14] F. Feng, X. Wang, R. Li, Cross-modal Retrieval with Correspondence Autoencoder, in: *ACM MM*, 7–16, 2014.
- [15] A. Habibian, T. Mensink, C. G. Snoek, VideoStory: A New Multimedia Embedding for Few-Example Recognition and Translation of Events, in: *ACM International Conference on Multimedia (MM)*, 17–26, 2014.
- [16] A. Eisenschat, L. Wolf, Linking Image and Text with 2-Way Nets, in: *IEEE CVPR*, 4601–4611, 2017.
- [17] J. Gu, J. Cai, S. Joty, L. Niu, G. Wang, Look, Imagine and Match: Improving Textual-Visual Cross-Modal Retrieval with Generative Models, in: *IEEE CVPR*, 2018.
- [18] D. He, Y. Xia, T. Qin, L. Wang, N. Yu, T. Liu, W.-Y. Ma, Dual Learning for Machine Translation, in: *NIPS*, 820–828, 2016.
- [19] J.-Y. Zhu, T. Park, P. Isola, A. A. Efros, Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks, in: *IEEE ICCV*, 2223–2232, 2017.
- [20] H. Hotelling, Relations between two sets of variates, *Biometrika* 28 (1936) 321–377.
- [21] P. L. Lai, C. Fyfe, Kernel and Nonlinear Canonical Correlation Analysis, *IEEE TPAMI* 10 (05) (2000) 365–377.
- [22] P. Mineiro, N. Karampatziakis, A Randomized Algorithm for CCA, in: *Neural Information Processing Systems (NIPS) workshop*, 2014.

- [23] T. Michaeli, W. Wang, , K. Livescu, Nonparametric Canonical Correlation Analysis, in: International Conference on Machine Learning (ICML), 1967–1976, 2016.
- [24] D. R. Hardoon, S. R. Szedmak, J. R. Shawe-taylor, Canonical Correlation Analysis: An Overview with Application to Learning Methods, *Neural Computation* 16 (12) (2004) 2639–2664.
- [25] Y. Gong, Q. Ke, M. Isard, S. Lazebnik, A Multi-View Embedding Space for Modeling Internet Images, Tags, and Their Semantics, *IJCV* 106 (2) (2014) 210–233.
- [26] V. Ranjan, N. Rasiwasia, C. V. Jawahar, Multi-Label Cross-Modal Retrieval, in: *IEEE ICCV*, 4094–4102, 2015.
- [27] A. Salvador, N. Hynes, Y. Aytar, J. Marin, F. Ofli, I. Weber, A. Torralba, Learning Cross-modal Embeddings for Cooking Recipes and Food Images, in: *IEEE CVPR*, 3020–3028, 2017.
- [28] L. Ma, Z. Chen, L. Xu, Y. Yan, Multimodal deep learning for solar radio burst classification, *Pattern Recognition* 61 (2017) 573–582.
- [29] S. Karaoglu, R. Tao, T. Gevers, A. W. M. Smeulders, Words Matter: Scene Text for Image Classification and Retrieval, *IEEE Transactions on Multimedia* 19 (5) (2017) 1063–1076.
- [30] X. Bai, M. Yang, P. Lyu, Y. Xu, Integrating Scene Text and Visual Appearance for Fine-Grained Image Classification with Convolutional Neural Networks, *CoRR* abs/1704.04613.
- [31] G. Andrew, R. Arora, K. Livescu, J. Bilmes, Deep Canonical Correlation Analysis, in: *ICML*, 1247–1255, 2013.
- [32] F. Yan, K. Mikolajczyk, Deep Correlation for Matching Images and Text, in: *IEEE CVPR*, 3441–3450, 2015.
- [33] B. Klein, G. Lev, G. Sadeh, L. Wolf, Associating Neural Word Embeddings With Deep Image Representations Using Fisher Vectors, in: *IEEE CVPR*, 4437–4446, 2015.
- [34] A. Karpathy, F.-F. Li, Deep Visual-Semantic Alignments for Generating Image Descriptions, in: *IEEE CVPR*, 3128–3137, 2015.
- [35] L. Ma, Z. Lu, L. Shang, H. Li, Multimodal Convolutional Neural Networks for Matching Image and Sentence, in: *IEEE ICCV*, 2623–2631, 2015.
- [36] R. Kiros, R. Salakhutdinov, R. S. Zemel, Unifying visual-semantic embeddings with multimodal neural language models, in: *NIPS workshop*, 2014.
- [37] Y. Liu, L. Liu, Y. Guo, M. S. Lew, Learning visual and textual representations for multimodal matching and classification, *Pattern Recognition* 84 (2018) 51–67.
- [38] Z. Zheng, L. Zheng, M. Garrett, Y. Yang, Y. Shen, Dual-Path Convolutional Image-Text Embedding, *CoRR* abs/1711.05535.
- [39] B. Wang, Y. Yang, X. Xu, A. Hanjalic, H. T. Shen, Adversarial Cross-Modal Retrieval, in: *ACM Multimedia*, 154–162, 2017.
- [40] S. Rastegar, M. Soleymani, H. R. Rabiee, S. Mohsen Shojaei, MDL-CW: A Multimodal Deep Learning Framework With Cross Weights, in: *IEEE CVPR*, 2601–2609, 2016.
- [41] V. Vukotić, C. Raymond, G. Gravier, Bidirectional Joint Representation Learning with Symmetrical Deep Neural Networks for Multimodal and Crossmodal Applications, in: *ICMR*, 343–346, 2016.
- [42] E. Kodirov, T. Xiang, S. Gong, Semantic Autoencoder for Zero-Shot Learning, in: *IEEE CVPR*, 3174–3183, 2017.
- [43] Y. Huang, Q. Wu, L. Wang, Learning Semantic Concepts and Order for Image and Sentence Matching, in: *IEEE CVPR*, 2018.
- [44] F. Wang, Q. Huang, L. Guibas, Image Co-Segmentation via Consistent Functional Maps, in: *IEEE ICCV*, 849–856, 2013.
- [45] T. Zhou, P. Krähenbühl, M. Aubry, Q. Huang, A. A. Efros, Learning Dense Correspondence via 3D-guided Cycle Consistency, in: *IEEE CVPR*, 117–126, 2016.

- [46] C. Godard, O. Mac Aodha, G. J. Brostow, Unsupervised Monocular Depth Estimation with Left-Right Consistency, in: IEEE CVPR, 270–279, 2017.
- [47] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative Adversarial Nets, in: NIPS, 2672–2680, 2014.
- [48] Z. Yi, H. Zhang, P. Tan, M. Gong, DualGAN: Unsupervised Dual Learning for Image-to-Image Translation, in: IEEE ICCV, 2849–2857, 2017.
- [49] T. Kim, M. Cha, H. Kim, J. K. Lee, J. Kim, Learning to Discover Cross-Domain Relations with Generative Adversarial Networks, in: ICML, 1857–1865, 2017.
- [50] R. Felix, B. V. Kumar, I. Reid, G. Carneiro, Multi-modal Cycle-Consistent Generalized Zero-Shot Learning, in: ECCV, 21–37, 2018.
- [51] G. Zhang, M. Kan, S. Shan, X. Chen, Generative Adversarial Network with Spatial Attention for Face Attribute Editing, in: V. Ferrari, M. Hebert, C. Sminchisescu, Y. Weiss (Eds.), ECCV, 422–437, 2018.
- [52] H. Tang, W. Wang, D. Xu, Y. Yan, N. Sebe, GestureGAN for Hand Gesture-to-Gesture Translation in the Wild, in: ACM Multimedia, 774–782, 2018.
- [53] X. Chen, C. L. Zitnick, Mind’s eye: A recurrent visual representation for image caption generation, in: IEEE CVPR, 2422–2431, 2015.
- [54] L. van der Maaten, G. Hinton, Visualizing Data Using t-SNE, *Journal of Machine Learning Research* 9 (2008) 2579–2605.
- [55] P. Young, A. Lai, M. Hodosh, J. Hockenmaier, From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions, in: ACL, 67–78, 2014.
- [56] K. Nandakumar, Y. Chen, S. C. Dass, A. Jain, Likelihood Ratio-Based Biometric Score Fusion, *IEEE TPAMI* 30 (2) (2008) 342–347.
- [57] L. Zheng, S. Wang, L. Tian, F. He, Z. Liu, Q. Tian, Query-Adaptive Late Fusion for Image Search and Person Re-Identification, in: IEEE CVPR, 1741–1750, 2015.
- [58] A. Bansal, S. Ma, D. Ramanan, Y. Sheikh, Recycle-GAN: Unsupervised Video Retargeting, in: ECCV, 2018.
- [59] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, T. Darrell, CyCADA: Cycle-Consistent Adversarial Domain Adaptation, in: ICML, 1989–1998, 2018.
- [60] J. Mao, W. Xu, Y. Yang, J. Wang, Z. Huang, A. Yuille, Deep Captioning with Multimodal Recurrent Neural Networks (m-RNN), in: ICLR, 2015.
- [61] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollar, C. L. Zitnick, Microsoft COCO: Common Objects in Context, in: ECCV, 740–755, 2014.
- [62] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, L. Fei-Fei, ImageNet Large Scale Visual Recognition Challenge, *IJCV* 115 (3) (2015) 211–252.
- [63] A. Conneau, D. Kiela, H. Schwenk, L. Barrault, A. Bordes, Supervised Learning of Universal Sentence Representations from Natural Language Inference Data, in: EMNLP, 670–680, 2017.
- [64] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, T. Darrell, Caffe: Convolutional Architecture for Fast Feature Embedding, in: ACM MM, 675–678, 2014.
- [65] X. Lin, D. Parikh, Leveraging Visual Question Answering for Image-Caption Ranking, in: ECCV, 261–277, 2016.
- [66] G. Lev, G. Sadeh, B. Klein, L. Wolf, RNN Fisher Vectors for Action Recognition and Image Annotation, in: ECCV, 833–850, 2016.
- [67] F. Faghri, D. J. Fleet, R. Kiros, S. Fidler, VSE++: Improved Visual-Semantic Embeddings, *CoRR* abs/1707.05612.
- [68] R. Kiros, Y. Zhu, R. R. Salakhutdinov, R. Zemel, R. Urtasun, A. Torralba, S. Fidler, Skip-Thought Vectors, in: NIPS, 3294–3302, 2015.

- [69] I. Vendrov, R. Kiros, S. Fidler, R. Urtasun, Order-embeddings of images and language, in: ICLR, 2016.
- [70] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed Representations of Words and Phrases and their Compositionality, in: NIPS, 3111–3119, 2013.