

Plug-and-Play Regulators for Image-Text Matching

Haiwen Diao, Ying Zhang, Wei Liu, Xiang Ruan and Huchuan Lu

Abstract—Exploiting fine-grained correspondence and visual-semantic alignments has shown great potential in image-text matching. Generally, recent approaches first employ a cross-modal attention unit to capture latent region-word interactions, and then integrate all the alignments to obtain the final similarity. However, most of them adopt one-time forward association or aggregation strategies with complex architectures or additional information, while ignoring the regulation ability of network feedback. In this paper, we develop two simple but quite effective regulators which efficiently encode the message output to automatically contextualize and aggregate cross-modal representations. Specifically, we propose (i) a Recurrent Correspondence Regulator (RCR) which facilitates the cross-modal attention unit progressively with adaptive attention factors to capture more flexible correspondence, and (ii) a Recurrent Aggregation Regulator (RAR) which adjusts the aggregation weights repeatedly to increasingly emphasize important alignments and dilute unimportant ones. Besides, it is interesting that RCR and RAR are “plug-and-play”: both of them can be incorporated into many frameworks based on cross-modal interaction to obtain significant benefits, and their cooperation achieves further improvements. Extensive experiments on MSCOCO and Flickr30K datasets validate that they can bring an impressive and consistent R@1 gain on multiple models, confirming the general effectiveness and generalization ability of the proposed methods.

Index Terms—Image-text matching, Recurrent correspondence regulator, Recurrent aggregation regulator, Cross-modal attention, Similarity aggregation, Plug-and-play operation.

I. INTRODUCTION

Exploiting the interactions between vision and language has attracted great interests in past decades, and various applications have sprouted to associate vision and text such as video-text retrieval [1]–[3], visual question answering [4], image captioning [5], visual grounding [6], and visual commonsense reasoning [7]. Among them, image-text matching involves the transmission and measurement of the cross-modal information, and provides great help for other tasks, making it become an important branch in the computer vision research area.

Great efforts have been made to accurately establish the relationship between visual and textual observations. Early works such as [4], [8]–[17] attempted to map the whole image and the full sentence into a joint embedding space, where the similarity between different modalities can be directly

This work was supported in part by the National Key R&D Program of China under Grant No.2018AAA0102001 and National Natural Science Foundation of China under grant No.62293542, U1903215 and the Fundamental Research Funds for the Central Universities No.DUT22ZD210.

H. Diao and H. Lu are with School of Information and Communication Engineering, Dalian University of Technology, Dalian, 116024, China. (Email: diaohw@mail.dlut.edu.cn; lhchuan@dlut.edu.cn). Y. Zhang and W. Liu are with Tencent Holdings Limited, Shenzhen, 518054, China. (Email: yinggzhang@tencent.com; wl2223@columbia.edu). X. Ruan is with Tiwaki Company Limited, Kusatsu, 5258577, Japan. (Email: ruanxiang@tiwaki.com).

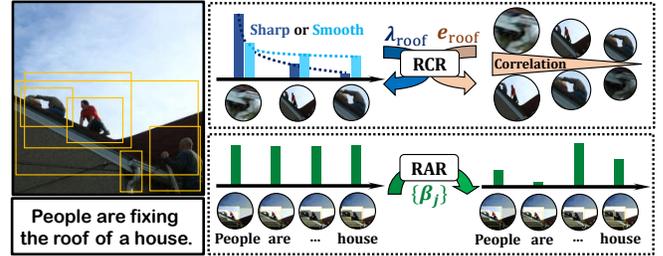


Fig. 1. Illustration of the proposed regulators. RCR progressively produces a plausible attention distribution between the word “roof” and its corresponding regions by adjusting the temperature λ and channel-wise factor e , while RAR constantly highlights significant alignments attended by each word and boosts itself step by step for more comprehensive aggregation.

measured. These approaches capture the global correspondence between an image and a sentence, while ignoring the importance of exploring fine-grained interactions across two modalities. To solve this problem, Lee *et al.* [18] proposed a cross-modal attention mechanism to explore region-word correspondences, which achieved impressive bi-directional retrieval results. Following it, many researchers are devoted to exploiting more accurate latent correspondence by either improving the cross-attention unit [19]–[21], or enhancing the cross-modal embeddings [22]–[26]. For example, Hu *et al.* [21] utilized a visual-semantic relation CNN to refine region-word interactions, while Chen *et al.* [23] reinforced the semantically related objects via encoding the image regions with a Recurrent Neural Network in order of the matching-word positions. Another thread of works focus more on the matching stage and infer the final similarity via aggregating all the alignments. Most existing approaches adopted the strategy of averaging all the cosine similarities between local alignments [18], [19], [21], [22], [26], which has achieved satisfying performance but remains far from being optimal. Chen *et al.* [23] introduced the self-attention mechanism to weight the image-word and object-text similarities separately, while [27], [28] performed graph reasoning on the matching vectors to gain more benefits over the simple cosine aggregation.

However, we observe that most approaches targeting cross-modal interactions focus on developing the interactive capability with one-time forward architectures [19], [21] or incorporating various information, such as intra-inter relationship [20], [25], [29], [30], object position [22], [27], semantic order [23], phrase structure [27], [31], while neglecting the regulation ability of network itself which learns from the message feedback and then leads to accurate and dynamic optimization. On the other hand, most of the above methods adopt the equivalent cross-modal alignment aggregation strategy for diverse region/word semantics and positive/negative pairs, lacking the ability to refine undesirable associations

and capture complicated matching patterns. In this paper, we introduce a regulator mechanism defined by [32]–[35] where the network can be improved by adaptively optimizing the forward learning process with plausible backward feedback loops, and validate that an elaborate regulation operation can make a vast difference in obtaining accurate interactions and conducting optimal aggregations across modalities requiring no additional data and complicated structures.

To be more specific, we propose a Recurrent Correspondence Regulator (RCR) and a Recurrent Aggregation Regulator (RAR) to progressively promote the image-text matching process, as shown in Fig. 1. The RCR learns adaptive attention factors for each specific word/region to refine the cross-modal attention unit iteratively, acquiring more plausible attention distributions for semantically diverse words/regions in various image-text pairs. The RAR starts with averaging all the alignments and then updates the aggregation weights guided by the aggregated alignment in the previous step, which increasingly emphasizes important alignments and gradually reduces the interference of unimportant ones to predict more precise similarity scores. An important and attractive property of the proposed RCR and RAR is “plug-and-play”: both of them can be seamlessly inserted into many existing methods based on cross-modal interaction to achieve remarkable improvements, and their cooperation brings greater benefits. Moreover, we experimentally verify that even with the simplest framework, the plug-ins of RCR and RAR enable the model [18] to achieve state-of-the-art results on MSCOCO and Flickr30K. In summary, our main contributions are three-fold:

- We propose a Recurrent Correspondence Regulator (RCR) to dynamically renew the cross-attention unit for better correspondence exploitation. It learns adaptive attention factors for each word/region to generate a more plausible attention distribution in accordance with its semantics and associated image-text pairs.
- We propose a Recurrent Aggregation Regulator (RAR) to repeatedly calibrate the weights for more discriminative similarity aggregation. It progressively reweights word/region-attended alignments directed by earlier guidance alignment to highlight more important alignments.
- The RCR and RAR can be applied to various approaches for image-text matching separately or jointly to achieve significant improvements, indicating the effectiveness and generalization ability of the proposed approach.

II. RELATED WORK

A. Cross-modal Attention

Cross-modal attention is first developed by Lee *et al.* [18] to discover all possible word-region alignments for image-text matching. With spectacular achievements, it attracts numerous researchers to make further explorations on enhancing cross-modal embeddings or improving attention units. Specifically, the former works attempt to facilitate word-region correspondence by enriching the instance features with region position [22], semantic order [23], scene graph [36], [37] and intra-inter correlation [20], while the latter methods directly develop more fine-grained interactions across modalities, such

as relation CNN [21], focal attention [19] and cross-graph attention [27], [37]. In particular, some works [27], [36], [37] introduced scene graphs with explicit attribute and relation information, and then constructed an inter-graph attention between graph nodes. Besides, Qu *et al.* [38] developed a routing mechanism to realize dynamic modality interaction, while Zhang *et al.* [39] exploited both the matched and mismatched effects for a comprehensive image-text relationship. Compared with previous works recurrently updating query features [26], instance fusion [4], [29], and context enhancement [20], the RCR directly adjusts the attention factors including the channel-wise weight vector and the softmax temperature, enabling attention weights to adapt to diverse semantic regions/words with different word/region sets from various image-text pairs. To be specific, for positive image-text pairs, the attention weights between each word/region and its corresponding regions/words are more precise and discriminative by the RCR, leading to a tighter distance between image and text. More importantly for negative pairs with completely irrelevant instances, the inner-product-like weights between the paired features adopted by the existing attention designs still emphasize the closest contents across modalities and capture the so-called “region-word correlations” in the latent space, inevitably increasing the image-text similarity and reducing the gap from the positive pairs. In contrast, the RCR progressively adjusts word-region relevance by the learned attention factors to generate appropriate attention distributions targeting diverse regions/words, meanwhile decoupling the attention weights from the final similarity measurement and producing larger gaps between matched and unmatched pairs.

B. Similarity Aggregation

Existing approaches [14], [15], [17], [40]–[42] based on mono-modal representation map the image and text features into a joint space and adopt the cosine distance as the measurement, while a great many methods [18], [19], [21], [22], [26] based on cross-modal interaction first obtain the pairwise features across modalities and then employ the average operation to fuse the cosine similarity between all the word-region alignments. Considering that various instances and hierarchical relevance have different importance in characterizing the cross-modality relations, Chen *et al.* [23] designed a self-attention module to integrate all cosine distances attended by regions or words, while Ji *et al.* [43] explored both fragment-wise and context-wise similarity scores to yield sufficient visual-semantic alignments between image and text. Besides for more powerful distance representations, some methods [27], [28], [44] introduced a vector-based similarity function and performed the matching pattern with graph reasoning, which have achieved great improvements at the cost of high complexity. Note that Liu *et al.* [27] not only needs to parse additional visual/textual graphs, but also fails to measure diverse alignments by aggregating them with average pooling operation. Moreover, Diao *et al.* [28] requires a high-quality holistic alignment to better guide the integration procedure of fragment-wise alignments. In contrast, our RAR employs a recurrent aggregation process without any extra

supplement and precondition, and validates that an iterative guidance alignment encoding early matching information can yields more appropriate weights and effectively facilitate the aggregation process for various alignments.

C. Plug-and-Play Methods

The modules that enable efficient integration into main frameworks are referred to as "plug-and-play" approaches. In recent years, the plug-and-play manners have attracted more attention in various fields, including image restoration [45]–[48], visual captioning [49], [50], visual question answering [51], [52], and video-text matching [53], [54]. By decoupling a specific problem from overall optimization objectives, they greatly simplify the integration process of each module, and improve flexibility and generalizability on new frameworks, thus accelerating the developments over other more sophisticated applications. It is worth noting that the method most relevant to us is GPO [17], which attempts to improve the mono-modal feature encoders and learn the best pooling strategy to integrate mono-modal instance features into a holistic embedding, while our regulators aim to generalize over various cross-modal interaction methods and promote multi-modal attention and similarity aggregation.

III. BACKGROUND

In this section, we briefly review the Stacked Cross Attention Network (SCAN) [18], which serves as the pioneer in exploiting word-region correspondences and alignments for image-text matching task. The whole architecture consists of four aspects: Feature Extraction, Cross-modal Attention, Similarity Computation, and Objective Function.

A. Feature Extraction

Image Representation. Given an image, the Faster R-CNN [55] model pretrained on Visual Genome [56] is first utilized to detect K salient regions with bottom-up attention [5], followed with a linear layer transforming each region feature into a d -dimensional vector. Therefore the image is encoded as a set of region features $\mathbf{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_K\}$, with $\mathbf{v}_i \in \mathbb{R}^d$ denoting the feature of i -th region.

Text Representation. Given a sentence with L words, we represent each word with a one-hot vector by random initialization, and map it into a 300-dimensional word embedding, followed by a bi-directional GRU [57] to integrate the bidirectional contextual information. The final text feature is computed by averaging the forward and backward hidden states to obtain $\mathbf{T} = \{\mathbf{t}_1, \dots, \mathbf{t}_L\}$, and $\mathbf{t}_j \in \mathbb{R}^d$ indicates the representation of j -th word.

B. Cross-modal Attention

Here, we only depict the text-to-image (T2I) attention in detail, and the image-to-text attention (I2T) performs similar operations. Given a set of region features $\mathbf{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_K\}$ and word features $\mathbf{T} = \{\mathbf{t}_1, \dots, \mathbf{t}_L\}$, the attention unit first computes the cosine similarities between all word-region pairs:

$$c_{ij}^{(0)} = R(\mathbf{v}_i, \mathbf{t}_j | \mathbb{1}^d) = \frac{\mathbf{v}_i^\top (\mathbb{1}^d \odot \mathbf{t}_j)}{\|\mathbf{v}_i\| \|\mathbf{t}_j\|}, \quad (1)$$

where $R(\cdot, \cdot | \mathbb{1}^d)$ indicates the cosine similarity function which computes the inner product weighted by the channel-wise all-ones vector $\mathbb{1}^d$. The attention weights are then calculated by a softmax function as

$$\alpha_{ij}^{(0)} = \frac{\exp(\lambda \bar{c}_{ij}^{(0)})}{\sum_{i=1}^K \exp(\lambda \bar{c}_{ij}^{(0)})}, \quad (2)$$

where $\bar{c}_{ij}^{(0)} = [c_{ij}^{(0)}]_+ / \sqrt{\sum_{j=1}^L [c_{ij}^{(0)}]_+^2}$ with $[x]_+ = \max(x, 0)$, and λ is the temperature of the softmax function. Here, $\alpha_{ij}^{(0)}$ is the normalized attention weight capturing the correspondence between the j -th word and its related regions, and thus the image feature attended by each word can be obtained via

$$\hat{\mathbf{v}}_j^{(0)} = \sum_{i=1}^K \alpha_{ij}^{(0)} \mathbf{v}_i. \quad (3)$$

For simplicity, we define the cross-modal attention unit as

$$\hat{\mathbf{v}}_j^{(0)} = \text{CMA}(\mathbf{t}_j, \mathbf{V} | \mathbb{1}^d, \lambda), \quad (4)$$

where the integrated image feature $\hat{\mathbf{v}}_j^{(0)}$ represents the related image regions with respect to j -th word under the fixed attention factors, including a channel-wise weight vector $\mathbb{1}^d$ and a softmax temperature λ .

C. Similarity Computation

The final image-text similarity is computed by averaging all the cosine similarities between $\hat{\mathbf{v}}_j^{(0)}$ and \mathbf{t}_j as

$$\mathcal{S}_{T2I} = \frac{1}{L} \sum_{j=1}^L R(\hat{\mathbf{v}}_j^{(0)}, \mathbf{t}_j | \mathbb{1}^d). \quad (5)$$

Similarly, the predicted similarity score by I2T attention is denoted as \mathcal{S}_{I2T} , and the combination of these two scores usually produces greater retrieval results.

D. Objective Function

Given a matched image-text pair (\mathbf{V}, \mathbf{T}) , the hard ranking loss [9] with online negative mining only takes account of the nearest negatives $(\tilde{\mathbf{T}}, \tilde{\mathbf{V}})$ within a mini-batch \mathcal{D} . The similarity of positive pairs should be higher than that of negative pairs by a fixed margin value γ , which is formulated as

$$\mathcal{L} = \sum_{(\mathbf{V}, \mathbf{T}) \in \mathcal{D}} [\gamma + \mathcal{S}(\mathbf{V}, \tilde{\mathbf{T}}) - \mathcal{S}(\mathbf{V}, \mathbf{T})]_+ + [\gamma + \mathcal{S}(\tilde{\mathbf{V}}, \mathbf{T}) - \mathcal{S}(\mathbf{V}, \mathbf{T})]_+, \quad (6)$$

where $\mathcal{S}(\cdot, \cdot)$ represents the matching score of an image-text pair computed by the aforementioned network.

IV. METHODOLOGY

In this section, we will elaborate on the proposed Recurrent Correspondence Regulator (RCR) and Recurrent Aggregation Regulator (RAR) based on the cross-modal attention unit from SCAN [18]. These two regulators can effectively explore the regulatory capacity of the network itself and in turn significantly facilitate the learning process by exploiting the well-designed alignment feedback. For simplicity, we take the T2I attention to describe the proposed regulation strategies, which can be applied to the I2T attention in the same way.

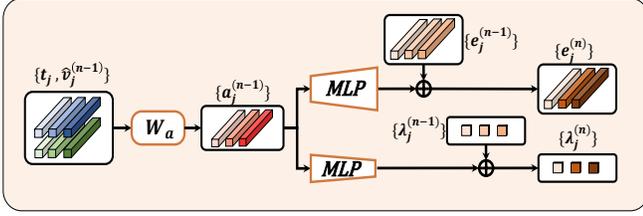


Fig. 2. Illustration of RCR that refines the cross-modal interactions via learning new channel-wise weight vectors and softmax temperature.

A. Recurrent Correspondence Regulator

The Recurrent Correspondence Regulator (RCR) learns specific attention factors for each word in a recurrent manner, aiming to refine the correspondence between each word and all the regions. In Fig. 2 with the word feature t_j and its related image feature $\hat{v}_j^{(0)}$, we first construct the alignment vector $\mathbf{a}_j^{(0)}$ following [58] with respect to t_j via

$$\mathbf{a}_j^{(0)} = \frac{\mathbf{W}_a |t_j - \hat{v}_j^{(0)}|^2}{\left\| \mathbf{W}_a |t_j - \hat{v}_j^{(0)}|^2 \right\|_2}, \quad (7)$$

where $\mathbf{W}_a \in \mathbb{R}^{m \times d}$ is a linear transformation, and $\mathbf{a}_j^{(0)} \in \mathbb{R}^m$ encodes the element-wise differences and fine-grained relationships between t_j and $\hat{v}_j^{(0)}$. With the comprehensive alignment encoding across two modalities, the alignment vector $\mathbf{a}_j^{(0)}$ is utilized to learn adaptive attention factors with multi-layer perceptron (MLP) for the next word-region interaction:

$$\begin{aligned} \mathbf{e}_j^{(1)} &= \left[\sigma(\mathbf{W}_{e'}(\sigma(\mathbf{W}_e \mathbf{a}_j^{(0)} + \mathbf{b}_e)) + \mathbf{b}_{e'}) + \mathbf{e}_j^{(0)} \right]_{-1}^{+1}, \\ \lambda_j^{(1)} &= \left[\mathbf{W}_{\lambda'}(\sigma(\mathbf{W}_\lambda \mathbf{a}_j^{(0)} + \mathbf{b}_\lambda)) + \mathbf{b}_{\lambda'} + \lambda_j^{(0)} \right]_+, \end{aligned} \quad (8)$$

where $\mathbf{W}_{\{\cdot\}}$ and $\mathbf{b}_{\{\cdot\}}$ are several learnable parameters, $\sigma(\cdot)$ indicates the tanh activation, and $[x]_{-1}^{+1}$ clips the value x to be within $[-1, +1]$. Note that each value of the vector $\mathbf{e}_j^{(1)}$ ranges from -1 to 1, reweighing the channel-wise negative or positive correlation between t_j and $\hat{v}_j^{(0)}$. Besides, the scalar $\lambda_j^{(1)}$ belongs to $[0, +\infty)$, controlling the word-wise smoothness or sharpness of attention distribution in relation to t_j .

Then we refine the word-region correspondence in the next step by separately reformulating Eq. (1)-(3) as

$$c_{ij}^{(1)} = R(\mathbf{v}_i, t_j | \mathbf{e}_j^{(1)}) = \frac{\mathbf{v}_i^\top (\mathbf{e}_j^{(1)} \odot t_j)}{\|\mathbf{v}_i\| \|t_j\|}, \quad (9)$$

$$\alpha_{ij}^{(1)} = \frac{\exp(\lambda_j^{(1)} \bar{c}_{ij}^{(1)})}{\sum_{i=1}^K \exp(\lambda_j^{(1)} \bar{c}_{ij}^{(1)})}, \quad (10)$$

$$\hat{v}_j^{(1)} = \sum_{i=1}^K \alpha_{ij}^{(1)} \mathbf{v}_i, \quad (11)$$

where $\mathbf{e}_j^{(1)} \in \mathbb{R}^d$ is the adaptive channel-wise weight vector learned to rectify the correlation between \mathbf{v}_i and t_j , which we term as $R(\cdot, \cdot | \mathbf{e}_j^{(1)})$, and $\lambda_j^{(1)} \in \mathbb{R}^1$ is the adaptive word-wise softmax temperature which adjusts the attention distribution.

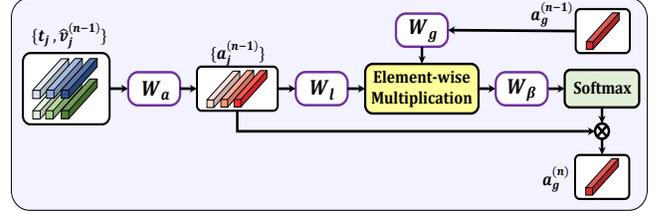


Fig. 3. Illustration of RAR that updates the aggregation weights under the guidance of the holistic alignment vector in the previous step.

$\hat{v}_j^{(1)}$ is the new integrated image feature with respect to j -th word. \odot denotes the element-wise multiplication.

The above equations illustrate how to update the word-attended image feature via learning two new attention factors in a single step. Similarly, the regulation process can be extended to multiple runs for further refinement. In this paper, we simplify the RCR as

$$\mathbf{e}_j^{(n)}, \lambda_j^{(n)} = \mathbf{RCR}(t_j, \hat{v}_j^{(n-1)}, \mathbf{e}_j^{(n-1)}, \lambda_j^{(n-1)}), \quad (12)$$

and plug it into the cross-modal attention unit via

$$\hat{v}_j^{(n)} = \mathbf{CMA}(t_j, \mathbf{V} | \mathbf{RCR}(t_j, \hat{v}_j^{(n-1)}, \mathbf{e}_j^{(n-1)}, \lambda_j^{(n-1)})), \quad (13)$$

where $\hat{v}_j^{(n)}$ is the updated image feature attended by j -th word in the n -th regulation step.

Discussion. For word-region interactions, **1)** most existing approaches compute the one-time forward procedures with the fixed and uniform factors, which obviously lack the regulation ability to adapt itself to various words with diverse semantics. In contrast, the RCR first generates the constructed alignment that records the abundant correlation between each word and all related regions from the previous step, which in turn reweighs the weight vector and temperature value concerning each word to refine the corresponding attention distribution. **2)** Early works are always inclined to align the words with potentially "closest" regions in the comparable space even for negative image-text pairs. We assume that the words from positive pairs should focus more on specific and relevant regions, while the ones from negative pairs should attend to "completely irrelevant" regions. From the above perspective, the RCR can dynamically update the channel-wise measure and refine the numerical value of word-region relevance, thus leading to larger gaps between matched and unmatched pairs and greater capability in modeling complex matching patterns.

B. Recurrent Aggregation Regulator

The Recurrent Aggregation Regulator (RAR) aggregates the word-region alignments in a recurrent manner by progressively optimizing the aggregation weights guided by the holistic alignment at the early step in Fig. 3. Given the word-attended alignment vector \mathbf{a}_j in Eq. (7), we initialize a guidance alignment with

$$\mathbf{a}_g^{(0)} = \frac{1}{L} \sum_{j=1}^L \mathbf{a}_j, \quad (14)$$

which actually performs the average pooling with the aggregation weight to be $1/L$ for each alignment. Instead of directly

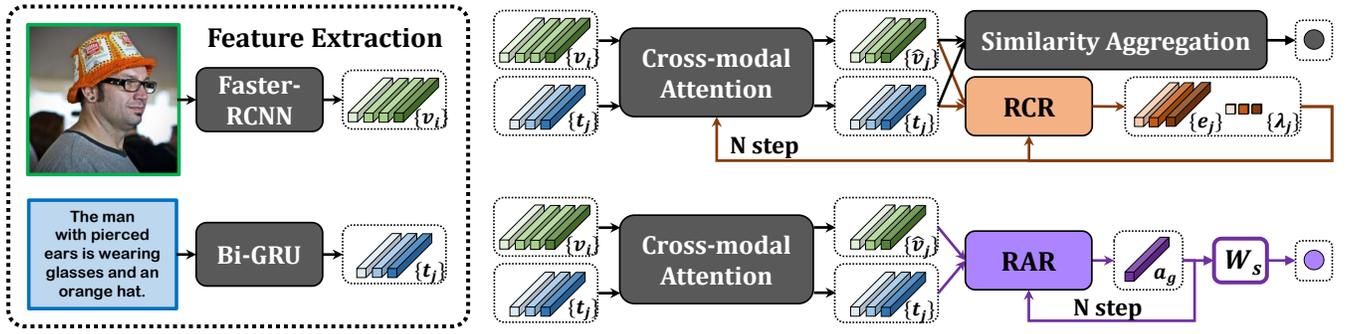


Fig. 4. Illustration of plug-and-play operation with our regulators. For independent application, the RCR facilitates region-word correspondence and preserves the raw similarity calculation, while the RAR promotes more accurate similarity prediction and retains the original cross-modal interaction.

using the averaged alignment $\mathbf{a}_g^{(0)}$ for inferring the similarity score, we iteratively update the aggregation weights under the guidance of $\mathbf{a}_g^{(0)}$ in the previous step:

$$\begin{aligned} \mathbf{u}_j^{(1)} &= \tanh(\mathbf{W}_g \mathbf{a}_g^{(0)}) \odot \tanh(\mathbf{W}_l \mathbf{a}_j), \\ \beta_j^{(1)} &= \frac{\exp(\mathbf{W}_\beta \mathbf{u}_j^{(1)})}{\sum_{i=1}^L \exp(\mathbf{W}_\beta \mathbf{u}_i^{(1)})}, \end{aligned} \quad (15)$$

where $\mathbf{W}_g \in \mathbb{R}^{m \times m}$, $\mathbf{W}_l \in \mathbb{R}^{m \times m}$, $\mathbf{W}_\beta \in \mathbb{R}^m$ are learnable parameters. The initial guidance alignment $\mathbf{a}_g^{(0)}$ in Eq. (14) is then updated as follows:

$$\mathbf{a}_g^{(1)} = \sum_{j=1}^L \beta_j^{(1)} \mathbf{a}_j, \quad (16)$$

where $\beta_j^{(1)}$ is the updated aggregation weight for the j -th alignment. For simplicity, we formulate the Recurrent Aggregation Regulator (RAR) as

$$\mathbf{a}_g^{(n)} = \text{RAR}(\mathbf{a}_g^{(n-1)}, \mathbf{A}), \quad (17)$$

with $\mathbf{A} = \{\mathbf{a}_1, \dots, \mathbf{a}_L\}$ indicating all the alignments constructed from the T2I attention as with Eq. (7).

The final similarity score can be inferred from $\mathbf{a}_g^{(n)}$ with a fully-connected layer as

$$\mathcal{S}_{T2I}^{\text{RAR}} = \text{sigmoid}(\mathbf{W}_s \mathbf{a}_g^{(n)}), \quad (18)$$

where $\mathbf{W}_s \in \mathbb{R}^m$ is a learnable parameter, and $\text{sigmoid}(\cdot)$ aims to output a similarity score within $[0, 1]$.

Discussion. Instead of averaging all the cosine similarities between all word features and attended image features as formulated in Eq. (5), the RAR goes one step further by iteratively aggregating the constructed alignments to recognize more comprehensive contents across modalities. Specifically, the RAR starts from the average aggregation, and in each regulation step it attempts to learn from the contextual message outputs from the previous step and balance the importance of each word-based alignment without no manual tuning. It is observed that the RAR increasingly emphasizes more on the alignments from more significant words, and gradually reduces the aggregation weights from unimportant ones. By this means, the network constantly adjusts the proportion of all the alignments and assigns more plausible aggregation weights, resulting in a more discriminative holistic alignment and more appropriate distance metrics in image-text matching.

C. Properties of RCR and RAR

Plug-and-Play on Multiple Models. The most attractive property of RCR and RAR is “plug-and-play”. To demonstrate their great applicability, we apply these two regulators to many existing methods based on cross-modal interaction:

Stacked Cross Attention (SCAN) [18] first computes all region-word similarities and aligns each region/word with its corresponding words/regions. The final similarity is obtained by averaging all region/word-based cosine distances.

Bidirectional Focal Attention (BFAN) [19] extends the generic attention by reassigning more fine-grained attention weight for each region-word pair and calculates the matching result by summing up region-based and word-based scores.

Position Focused Attention (PFAN) [22] enhances region features by introducing extra position information to promote region-word correspondences and integrates all region/word-attended cosine similarities as the prediction.

Cross-Modal Adaptive Message Passing (CAMP) [30] explores a region-word affinity matrix via inner product and transfers cross-modality contents to improve the region and word representations, which are then aggregated as the holistic image and text features to compute the final similarity.

Similarity Graph Reasoning and Attention Filtration (SGRAF) [28] adopts cosine similarities multiplied with a fixed temperature as region-word attention weights, followed by the complex graph and attention modules to map hierarchical similarity features into a matching score.

Fig. 4 illustrates how we plug the RCR or RAR into the above matching approaches. Specifically, cross-modal attention utilizes the cosine metric or inner product as region-word affinity weights, and outputs each region/word along with its related words/regions. With these paired features, the RCR first constructs the alignment vectors and then learns the corresponding weight vectors and temperature factors via Eq. (7)-(8), which in turn refine the region-word feature distances and optimize the cross-modal interaction via Eq. (9)-(11). Besides with a set of alignment vectors, the RAR progressively generates more appropriate weights between a guidance vector and all alignment vectors via Eq. (14)-(16), and facilitates more rational similarity aggregation processing. It turns out that such simple message feedback brings remarkable improvements on many cross-modal interaction works, and even achieves superior performance than related complicated counterparts.

TABLE I

RETRIEVAL RESULTS IN CHRONOLOGICAL ORDER. THE BEST TWO RESULTS ARE MARKED IN **BOLD** AND UNDERLINE. * ADOPTS WARM-UP STRATEGY AND TEXT SIZE AUGMENTATION, WHILE * DENOTES ENSEMBLE MODELS WITH HIGH RESOLUTION OF THE INPUT IMAGES.

Methods	Flickr30K 1K Test						MSCOCO 5-fold 1K Test						MSCOCO 5K Test					
	Sentence Retrieval			Image Retrieval			Sentence Retrieval			Image Retrieval			Sentence Retrieval			Image Retrieval		
	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
Faster-RCNN (ResNet-101 BUTD [5]) + Random Word Embedding Initialization																		
SCAN [18] _{ECCV18}	67.4	90.3	95.8	48.6	77.7	85.2	72.7	94.8	98.4	58.8	88.4	94.8	50.4	82.2	90.0	38.6	69.3	80.4
VSRN [14] _{JCCV19}	71.3	90.6	96.0	54.7	81.8	88.2	76.2	94.8	98.2	62.8	89.7	95.1	53.0	81.1	89.4	40.5	70.6	81.1
CAAN [25] _{CVPR20}	70.1	91.6	97.2	52.8	79.0	87.9	75.5	95.4	98.5	61.3	89.7	95.2	52.5	83.3	90.9	41.2	70.3	82.9
IMRAM [26] _{CVPR20}	74.1	93.0	96.6	53.9	79.4	87.2	76.7	95.6	98.5	61.7	89.1	95.0	53.7	83.2	91.0	39.7	69.1	79.8
MMCA [20] _{CVPR20}	74.2	92.8	96.4	54.8	81.4	87.8	74.8	95.6	97.7	61.6	89.8	95.2	54.0	82.5	90.7	38.7	69.7	80.8
GSMN [27] _{CVPR20}	76.4	<u>94.3</u>	97.3	57.4	82.3	<u>89.0</u>	78.4	<u>96.4</u>	<u>98.6</u>	<u>63.3</u>	90.1	95.7	-	-	-	-	-	-
SGRAF [28] _{AAAI21}	77.8	94.1	97.4	<u>58.5</u>	<u>83.0</u>	88.8	79.6	96.2	98.5	63.2	90.7	96.1	<u>57.8</u>	<u>84.9</u>	<u>91.6</u>	<u>41.9</u>	<u>70.7</u>	81.3
SHAN [43] _{JCAI21}	74.6	93.5	96.9	55.3	81.3	88.4	76.8	96.3	98.7	62.6	89.6	<u>95.8</u>	-	-	-	-	-	-
WCGI [40] _{JCCV21}	74.8	93.3	96.8	54.8	80.6	87.5	75.4	95.5	<u>98.6</u>	60.8	89.3	95.3	-	-	-	-	-	-
RCAR([18]_T2I)	<u>77.8</u>	93.6	96.9	57.2	82.8	88.5	78.2	96.3	98.4	62.2	89.6	95.3	57.4	83.8	91.0	40.7	69.8	80.4
RCAR([18]_I2T)	74.7	93.0	97.1	54.6	80.5	87.0	78.5	95.9	98.5	61.2	89.0	95.2	56.6	83.3	91.2	39.1	68.7	79.4
RCAR([18]_All)	78.7	94.6	97.6	59.5	84.0	89.5	80.6	96.6	<u>98.6</u>	64.1	<u>90.5</u>	<u>95.8</u>	59.6	85.8	92.4	42.5	71.7	<u>81.8</u>
SCAN [18] _{ECCV18*}	72.2	92.4	96.5	53.6	81.2	88.6	72.8	94.3	98.0	57.5	87.8	94.5	50.1	79.5	88.1	36.5	66.7	78.1
GPO [17] _{CVPR21*}	78.0	94.6	97.8	58.3	84.6	90.8	78.4	96.2	98.7	62.7	90.6	95.9	56.8	84.5	91.4	40.3	70.7	81.7
GPO [17] _{CVPR21**}	80.7	96.4	<u>98.3</u>	<u>60.8</u>	86.3	92.3	80.0	97.0	99.0	<u>64.8</u>	91.6	96.5	<u>59.8</u>	86.1	92.8	<u>42.7</u>	<u>72.8</u>	83.3
RCAR([18]_T2I*)	79.7	95.0	97.4	60.9	84.4	90.1	79.1	96.5	98.8	63.9	90.7	95.9	59.1	84.8	91.8	42.8	71.5	81.9
RCAR([18]_I2T*)	76.9	95.5	98.0	58.8	83.9	89.3	79.3	96.5	98.8	63.8	90.4	95.8	58.4	84.6	91.9	41.7	71.4	81.7
RCAR([18]_All*)	82.3	<u>96.0</u>	98.4	62.6	<u>85.8</u>	<u>91.1</u>	80.9	<u>96.9</u>	<u>98.9</u>	65.7	<u>91.4</u>	<u>96.4</u>	61.3	86.1	<u>92.6</u>	44.3	73.2	<u>83.2</u>

Algorithm 1 Cooperation of RCR and RAR (RCAR)

Input: Image features V , text features T , initial temperature λ and weight vector $\mathbb{1}^d$, and regulation steps N ;

Output: Final similarity score \mathcal{S}^{RCAR} ;

- 1: Compute $\hat{v}_j^{(0)}, j = 1, \dots, L$ with Eq. (4);
- 2: Compute $\mathbf{a}_j^{(0)}, j = 1, \dots, L$ with Eq. (7);
- 3: Compute $\mathbf{a}_g^{(0)}$ with Eq. (14);
- 4: **for** $n = 1$ to N **do**
- 5: Update $\hat{v}_j^{(n)}, j = 1, \dots, L$ with Eq. (13);
- 6: Update $\mathbf{a}_j^{(n)}, j = 1, \dots, L$ with Eq. (7);
- 7: Update $\mathbf{a}_g^{(n)}$ with Eq. (17);
- 8: **end for**
- 9: Compute \mathcal{S}^{RCAR} with Eq. (18);
- 10: **return** \mathcal{S}^{RCAR}

Cooperation of RCR and RAR. The RCR and RAR can cooperate with each other where the RCR is responsible for adjusting the cross-modal interaction and the RAR refines the alignment aggregation to achieve further improvements. In Algorithm 1, we introduce an easy combination as RCAR that performs these two regulations one-by-one. Note that their cooperation is pretty flexible, and more variants with experimental results can be found in Sec. V-C.

V. EXPERIMENTS

In this section, we first describe the detailed implementations and training settings, and then validate the great performance and generalization ability of two regulators.

A. Datasets and Settings

Datasets and Protocols. We utilize MSCOCO [59] and Flickr30K [60] that separately consist of 31,783 and 123,287 images, with each one annotated with 5 text descriptions. For Flickr30K, we split the dataset into 1,000 images for

validation, 1,000 images for testing, and the rest for training. For MSCOCO, we utilize 113,287 images for training, 5,000 images for validation, and 5,000 images for testing. We report the results by averaging over 5 folds of 1K test images and testing on the full 5K test images, respectively. In terms of evaluation metric, we measure the performance by the Recall@ \hat{K} ($R@K$) which measures the fraction of queries whose ground-truth is ranked among the closest \hat{K} results.

Implementation Details. The bottom-up detector [5] is used to generate the top $K=36$ region proposals with 2048 dimensions. Besides, we set the dimensions of word embedding, hidden state of BiGRU, and alignment vector as 300, $d=1024$, and $m=256$ respectively. The initial $\lambda^{(0)}=10$ and $e^{(0)}=\mathbb{1}^d$ are updated by two MLPs of Input(256)-FC(128)-Tanh-FC(1) and Input(256)-FC(512)-Tanh-FC(1024)-Tanh. The network is trained by the Adam optimizer [61] with a mini-batch size of 128. For MSCOCO, we set the learning rate to be 0.0002 for the first 10 epochs and 0.00002 for the next 10 epochs. For Flickr30K, the learning rate is set to be 0.0002 for 30 epochs and decayed by 0.1 for the next 10 epochs.

B. Quantitative Results and Analysis

We present the results with $N=2$ RCAR (i.e. 2-step RAR and 1-step RCR) with the simplest SCAN and improved SCAN* that adopts a warm-up strategy and text size augmentation as with [17]. We report the ensemble results of T2I and I2T models by averaging the individual scores offline.

Results on Flickr30K. TABLE I shows the retrieval results on Flickr30K. Compared with SCAN [18], our regulators can improve the absolute R@1 boost of 11.3% and 10.9% on sentence and image retrieval. Besides, the RCAR with the improved SCAN [18]* yields the bidirectional R@1 of 82.3% and 62.6% separately, and exceeds the best competitor GPO [17] by 4.3% and 4.3% under the same settings, indicating the significance of exploiting the regulation capabilities with adaptive correspondence and recurrent aggregation.

TABLE II

RETRIEVAL RESULTS OF PLUG-AND-PLAY RAR, RCR AND RCAR ON FLICKR30K WITH THE OFFICIAL CODES OF MULTIPLE APPROACHES.

Methods	#RAR	#RCR	Sen. Ret.		Ima. Ret.		Mem. (G)	Tim. (us)
			R@1	R@5	R@1	R@5		
BFAN [19]	X	X	68.1	91.4	50.8	78.4	11.2	11.8
	3	X	74.5	92.9	54.5	80.6	12.4	14.4
SGRAF [28]	X	X	77.8	94.1	58.5	83.0	12.4	9.4
	X	2	79.2	94.3	59.7	83.1	13.4	28.4
CAMP [30]	X	X	68.1	89.7	51.5	77.1	18.8	36.8
	X	2	75.1	93.2	56.0	81.3	20.1	78.5
	3	X	74.4	91.9	53.7	80.0	19.0	43.1
	2	1	76.3	93.3	57.1	81.6	19.4	55.3
PFAN [22]	X	X	70.0	91.8	50.4	78.7	10.4	8.8
	X	2	73.6	92.5	54.3	81.1	11.7	43.0
	3	X	78.1	94.1	58.4	82.9	10.8	12.6
	2	1	80.1	95.7	59.9	84.4	10.9	26.1

Results on MSCOCO. In TABLE I with 5-fold 1K test images, our RCAR can produce the state-of-the-art performance based on the simplest SCAN [18], and outweigh the SGRAF [28] by 1.0% and 0.9% on the most concerned R@1. Under the fair comparison, the improved version consistently surpasses the previous best method GPO [17] by 2.5% and 3.0% R@1 increases at two directions. With the larger and more compelling 5K test images, our RCAR with [18] and [18]* can further outperform the SGRAF [28] and GPO [17] by 1.8/0.6% and 4.5/4.0% R@1 improvements respectively, validating the superior performance and generalization capability in handling more complex matching patterns.

Plug-and-Play on Multiple Models. We attempt to apply our regulators on a series of representative works including BFAN [19], SGRAF [28], CAMP [30], and PFAN [22] on Flickr30K in TABLE II. **1)** Since BFAN [19] designs specific cross-modal attention which explores a novel bidirectional focal attention to eliminate irrelevant fragments from the shared semantic, we just plug 3-step RAR into the network and obtain the R@1 gains of 6.4% and 3.7% on BFAN. Note that when both applied with the RAR, the cross-modal attention unit from SCAN [18] refined by the RCR achieves much better performance (R@1=78.7/59.5%) than BFAN (R@1=74.5/54.5%), which further verifies the superior cross-modal correspondence by exploiting the regulation abilities of the network itself. **2)** SGRAF [28] employs graph reasoning and attention filtration to refine the cross-modal representations, but ignores the ability of cross-modal attention unit, which can work with the RCR to empower flexible region-word interactions. With 2-step RCR, SGRAF obtains a maximum 1.4% increase on R@1, reflecting general effectiveness with the complicated network. **3)** CAMP [30] and PFAN [22] integrate the cross-modal message flow and valuable position embedding separately to enhance the multi-modal representations, which can possess more powerful cross-modal interaction and aggregation actuated by our regulators. When the RCR/RAR/RCAR is introduced, the bidirectional R@1 can rise by 7.0/6.3/8.2% and 4.5/2.2/5.6% on CAMP, as well as 3.6/8.1/10.1% and 3.9/8.0/9.5% on PFAN on sentence and image retrieval respectively, demonstrating the strong compatibility and flexibility of our approach. **4) Computational cost.** Here, we report the memory and time consumption

TABLE III

RESIDUAL DESIGN OF THE REGULATORS WITH T2I ATTENTION ON FLICKR30K. RES DENOTES WHETHER TO USE RESIDUAL CONNECTION.

Model	Res	Step	Sentence Retrieval			Image Retrieval		
			R@1	R@5	R@10	R@1	R@5	R@10
Baseline	X	X	64.8	89.9	94.5	46.9	76.0	84.5
	X	1	69.9	91.1	95.4	51.9	79.9	86.0
RCR	X	2	61.4	85.1	91.2	43.0	72.5	81.7
	✓	1	73.0	93.3	97.4	55.3	80.2	86.9
	✓	2	74.3	93.3	97.1	56.6	81.3	87.8
RAR	X	2	75.7	93.6	97.2	56.0	81.8	88.0
	X	3	76.2	93.8	96.8	56.7	81.8	88.2
	✓	2	76.6	93.4	96.5	56.8	81.2	86.9
	✓	3	75.8	92.9	96.7	56.6	80.6	85.8

for prediction and average the additional cost for ensemble models. With 3-step RAR, the extra time increase of each image-text pair for BFAN/CAMP/PFAN is 2.6/6.3/3.8 *us* with the memory increase of 1.2/0.2/0.4 *G*, while with 2-step RCR, the extra cost for SGRAF/CAMP/PFAN of 19/41.7/34.2 *us* and 1.0/1.3/1.3 *G*. Besides, our RCAR brings the time and memory cost for CAMP/PFAN of 18.5/17.3 *us* and 0.6/0.5 *G*, and gains a good balance between accuracy and complexity.

C. Ablation Studies

In this section, we first report the configurations of our proposed regulators, as well as the initialization and optimization of the attention factors. Then, we delve into the RAR and RCR to display how the aggregation weights and cross-attention distributions are progressively refined. Finally, we also explore alternative strategies and architectures. All comparisons are implemented based on SCAN [18] unless otherwise noted.

Residual mechanism of the regulators. In TABLE III, we carry out critical analyses of the influence of residual architectures. The Baseline employs the T2I attention from SCAN [18] and averages all the cosine similarities as the final score. **1) Correspondence regulator.** Eq. (8) indicates that the current adaptive weight vector $e_j^{(n)}$ and softmax temperature $\lambda_j^{(n)}$ require the $e_j^{(n-1)}$ and $\lambda_j^{(n-1)}$ at the last step. Here, we remove these two variables to construct a no-residual version of the RCR. Compared with the residual structure, the RCR without residual design results in an obvious R@1 drop in TABLE III, indicating that RCR is inclined to predict offsets against the current state to adjust previous regulation dynamically. To be specific, 1-step RCR without residual fashion produces better results than Baseline. This is because in the beginning, each word shares the same initialization of a weight vector $e^{(0)}=\mathbb{1}^d$ and temperature $\lambda^{(0)}=10$, and the RCR barely infers the absolute value of these attention factors in the next step. However, after a 1-step adjustment, all the word-region interactions start from very distinct conditions (aligned or not) and attention states with regard to the particular words, making it difficult to further forecast absolute valuations. Therefore, the RCR with residual mechanism can better adjust the dynamic learning process and reduce the burden of one-time total optimization in a progressive manner. **2) Aggregation regulator.** Eq. (16) denotes that the current guidance alignment $a_g^{(n)}$ requires no need for the $a_g^{(n-1)}$ in the last iteration. Similarly, we average the early and

TABLE IV

IMPACT OF #RCR AND #RAR WITH T2I ATTENTION ON MSCOCO5K. LIMITED BY THE MACHINE, WE SET THE MAXIMUM STEP OF RCR TO 4.

#RAR	#RCR	Sentence Retrieval			Image Retrieval		
		R@1	R@5	R@10	R@1	R@5	R@10
\times	\times	44.0	75.5	85.5	32.6	62.0	74.4
1	\times	56.3	82.9	89.8	39.6	69.2	79.8
2	\times	56.2	83.1	90.7	40.0	69.4	79.9
3	\times	56.6	83.5	90.8	40.5	69.4	80.4
4	\times	56.6	83.3	90.9	40.4	69.4	80.2
\times	1	47.8	79.4	89.2	35.2	66.8	78.4
\times	2	53.4	81.9	90.4	38.4	68.7	79.9
\times	3	52.4	81.8	90.6	38.3	69.0	80.0
\times	4	53.1	83.3	90.9	38.3	69.3	79.8
2	1	57.4	83.8	91.0	40.7	69.8	80.4
3	2	56.8	83.8	91.0	40.8	70.0	80.5

TABLE V

IMPACT OF #RCR AND #RAR WITH I2T ATTENTION ON MSCOCO5K. LIMITED BY THE MACHINE, WE SET THE MAXIMUM STEP OF RCR TO 3.

#RAR	#RCR	Sentence Retrieval			Image Retrieval		
		R@1	R@5	R@10	R@1	R@5	R@10
\times	\times	43.4	74.8	84.8	32.0	61.8	74.2
1	\times	52.7	81.6	90.2	36.8	68.0	78.8
2	\times	53.3	82.1	90.5	38.1	68.4	79.2
3	\times	53.1	81.8	90.4	38.4	68.6	79.3
\times	1	48.5	80.5	89.4	34.9	66.7	78.5
\times	2	54.2	82.0	90.3	38.3	67.8	78.9
\times	3	56.8	83.2	91.0	39.4	68.8	79.3
2	1	56.6	83.3	91.2	39.1	68.7	79.4
3	2	56.0	83.2	91.3	39.3	69.1	79.4

learned alignments as the current guidance vector to build a residual version of the RAR. The RAR aims to construct better bootstrap guidance and assign appropriate aggregation weights with the original word-based alignments throughout the process. Therefore, we can discover that the RAR with residual structure fails to bring significant improvements with the same word-attended alignments at each step.

Hyperparameter tuning of #RAR and #RCR. TABLE IV, V, VI and VII demonstrate the evaluation results of different steps about our regulators. We establish the baseline without any regulator that only utilizes the cross-modal attention [18] and predicts the final score by averaging all the cosine distances via Eq. (5). **1) Aggregation regulator.** The RAR holds the original cross-modal attention unit and calculates the similarity by Eq. (18) with the alignments constructed from Eq. (7). For MSCOCO 5K test set, the RAR can steadily improve the R@1 on sentence and image retrieval by at most 12.6% and 7.9% based on T2I attention, as well as 9.9% and 6.4% based on I2T attention. For Flickr30K 1k test set, it can also boost the bidirectional R@1 with consistent gains of over 9.8/7.5% and 2.4/7.7% upon T2I and I2T attention, respectively. We can see that the RAR can generate more accurate and plausible image-text similarity measurements. **2) Correspondence regulator.** The RCR renews the region-word interactions to update aggregated features targeting the cross-modality instances iteratively, and keeps the raw prediction process through Eq. (5). Compared with the foundation models, the RCR can obtain the steady R@1 increases by maximum 9.4/5.8% (T2I) and 13.4/7.4% (I2T) on MSCOCO5K, and meanwhile 10.4/9.9% (T2I) and

TABLE VI

IMPACT OF #RCR AND #RAR WITH T2I ATTENTION ON FLICKR30K. LIMITED BY THE MACHINE, WE SET THE MAXIMUM STEP OF RCR TO 4.

#RAR	#RCR	Sentence Retrieval			Image Retrieval		
		R@1	R@5	R@10	R@1	R@5	R@10
\times	\times	64.8	89.9	94.5	46.9	76.0	84.5
1	\times	74.6	92.6	96.1	54.4	80.4	87.2
2	\times	75.7	93.6	97.2	56.0	81.8	88.0
3	\times	76.2	93.8	96.8	56.7	81.8	88.2
4	\times	74.8	92.7	96.8	56.3	81.5	86.6
\times	1	73.0	93.3	97.4	55.3	80.2	86.9
\times	2	74.3	93.3	97.1	56.6	81.3	87.8
\times	3	73.9	93.1	96.4	56.0	81.4	87.5
\times	4	75.2	94.1	97.8	56.8	81.8	87.8
2	1	77.8	93.6	96.9	57.2	82.8	88.5
3	2	76.8	94.1	97.1	57.1	83.0	88.2

TABLE VII

IMPACT OF #RCR AND #RAR WITH I2T ATTENTION ON FLICKR30K. LIMITED BY THE MACHINE, WE SET THE MAXIMUM STEP OF RCR TO 3.

#RAR	#RCR	Sentence Retrieval			Image Retrieval		
		R@1	R@5	R@10	R@1	R@5	R@10
\times	\times	66.7	89.1	94.0	41.2	72.7	81.8
1	\times	69.1	91.2	95.6	48.9	77.1	85.1
2	\times	70.8	90.3	95.7	51.9	78.5	85.9
3	\times	70.6	91.2	95.0	50.6	77.8	85.3
\times	1	69.0	93.3	96.8	53.0	80.4	86.2
\times	2	73.3	92.3	96.7	54.3	80.5	87.3
\times	3	75.7	93.0	97.4	56.8	81.6	88.1
2	1	74.7	93.0	97.1	54.6	80.5	87.0
3	2	74.2	92.4	96.5	54.7	80.6	86.9

9.0/15.6% (I2T) on Flickr30K, verifying that RCR is capable of exploiting more fine-grained and appropriate word-region associations. **3) Cooperative regulators.** We employ the one-by-one combination of RAR and RCR as described in Algorithm 1, indicating that the former always takes one more step than the latter. Compared with 2-step RAR and 1-step RCR, $N=2$ RCAR can further promote the R@1 at two directions by a large margin, demonstrating the good compatibility between RCR and RAR. Actually, their cooperations are pretty flexible. To take T2I attention on Flickr30K in TABLE VI as an example, an alternative strategy is to first perform 2-step RCR followed by 3-step RAR, which yields the competitive 77.5 and 57.8% R@1 (against 77.8 and 57.2% R@1 according to Algorithm 1) on sentence and image retrieval separately. Besides, we can also observe that larger #RCR and #RAR are not necessarily better, which may be due to the recurrent structure where a certain number of steps can saturate the performance of the network. **4) Computational cost.** The model size of single SCAN [18] is 12.2M, and each step of RAR, RCR, or RCAR brings the extra parameters of nearly 0.13M, 0.95M, or 1.12M. Using NVIDIA GeForce RTX 3090, the inference time of T2I-SCAN is 3.05 us for each image-text pair with an extra cost of 0.51/4.72/5.91 us per step by RAR/RCR/RCAR, while the predicted time of I2T-SCAN is 4.25 us with an additional cost of 1.19/8.21/10.54 us. From these experiments, we suggest step=2-3 for independent application and step=2(RAR)+1(RCR) for their cooperation, as they achieve a better trade-off between accuracy and complexity, and have proved general effectiveness and broad applicability on multiple state-of-the-art approaches in TABLE II.

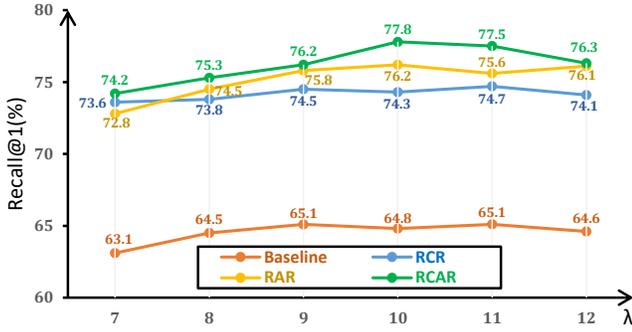


Fig. 5. Impact of the λ initialization with T2I attention on Flickr30K. The step of RAR, RCR and RCAR is set as 3, 2, and 2, respectively.

TABLE VIII

IMPACT OF THE ATTENTION FACTORS (λ , e) OPTIMIZATION WITH T2I ATTENTION ON FLICKR30K. WE ADOPT $N=2$ RCAR AS A REFERENCE.

Model	Sentence Retrieval			Image Retrieval		
	R@1	R@5	R@10	R@1	R@5	R@10
Fixed	75.7	93.6	97.2	56.0	81.8	88.0
Learnable	76.3	94.1	96.8	56.9	81.3	87.3
MLP(t_j)	75.9	93.3	96.5	56.6	81.2	87.3
MLP(a_j)	77.8	93.6	96.9	57.2	82.8	88.5

Initialization of the attention factors. Fig. 5 depicts the Recall@1(%) at sentence retrieval with varying λ value on Flickr30K. The λ determines the initial distributions of word-region interactions where a large one tends to retain only the highly correlated instances while a small one results in the interference from irrelevant instances. Here, we take T2I attention as an example and compare the performance variation with λ ranging from 7 to 12. We can see that our regulators can obtain the maximum performance benefit when $\lambda=10$, and achieve a consistent improvement among various settings, confirming the robustness and stability of our proposed method. It is worth noting that for simplicity and fair comparison, we directly set $e^{(0)}=\mathbb{1}^d$ and use exclusive λ of each work as $\lambda^{(0)}$ in TABLE II, which attempts to maintain the appropriate initialization of the incipient cross-attention unit.

Optimization of the attention factors. We investigate the different update strategies with $N=2$ RCAR in TABLE VIII. **1) Fixed:** We set the weight vector $e=\mathbb{1}^d$ and the softmax temperature $\lambda=10$ in the whole process; **2) Learnable:** The parameters e and λ are learnable during the training, with initialization of $\mathbb{1}^d$ and 10 in the beginning; **3) MLP(t_j):** The attention factors of each word are learned with the original word features; **4) MLP(a_j):** The attention factors of each word are learned with the constructed alignment vector. Note that **Learnable** achieves slightly better performance than **Fixed**, and adjusts the $\lambda = 11.23$ for maximum performance benefit in experiments. However, a common problem of the two methods is the lack of capability to refine parameters adaptively to handle the diversity of different words. **MLP(t_j)** produces even worse results than **Learnable**, which may be due to the lack of word-image alignment information, leading to difficulties in learning reasonable attention factors. In comparison, **MLP(a_j)** achieves the best performance by learning adaptive factors for each word, indicating the significance of exploiting the interaction feedback for better regulation.

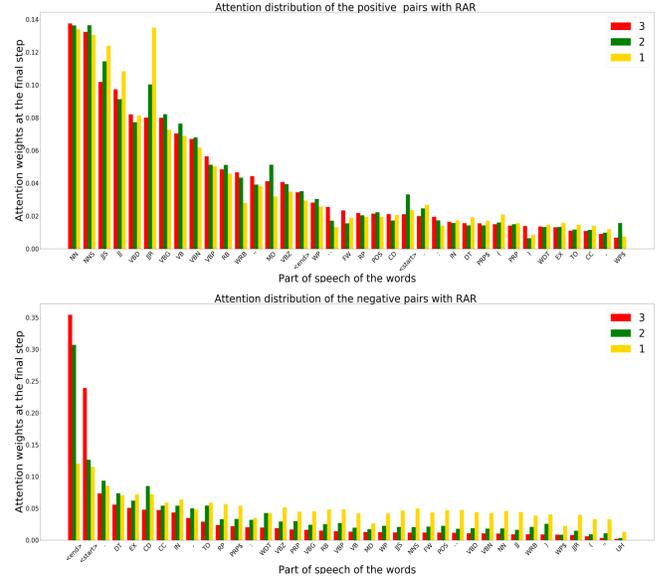


Fig. 6. Quantitative weight statistics of word-attended alignments by n-step RAR with T2I attention on Flickr30K. 1, 2, 3 indicate the steps of the RAR. The top and bottom represent positive and negative image-text pairs.

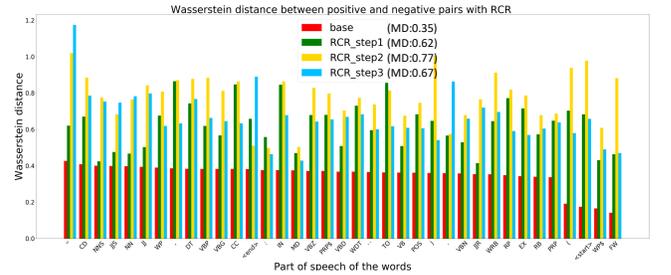


Fig. 7. Quantitative distance statistics of word-attended cosine similarities by n-step RCR with T2I attention on Flickr30K. Base denotes T2I-SCAN [18]. MD denotes the mean wasserstein distance with respect to all parts of speech.

Quantitative statistics of the regulators. Fig. 6 compares the attention weight distribution of different word-attended alignments with n-step RAR, while Fig. 7 displays the wasserstein distance of word-attended cosine similarities between positive and negative pairs with n-step RCR. Considering a very large range of vocabulary, we adopt NLTK toolkit and conduct statistical analyses with respect to the part of speech. **1) Aggregation regulator.** For the matched image-text pairs, the RAR attempts to reduce the interference of less-meaningful alignments and highlight the important ones attended by nouns, adjectives, and verbs that contain rich semantic information. On the other hand, "`<start>`" and "`<end>`" encode the global textual contextual representations by BiGRU, and their corresponding alignments are emphasized gradually for the unmatched pairs which reflect the holistic differences across modalities. **2) Correspondence regulator.** From Fig. 7, the mean distance learned by raw attention unit (T2I-SCAN [18]) is nearly 0.35, and the one by the RCR is over 0.6 (step=2 best) with all parts of speech between positive and negative image-text pairs. We assume that attention weights are computed by one-step forward interactions with the fixed weight vector and softmax temperature, which

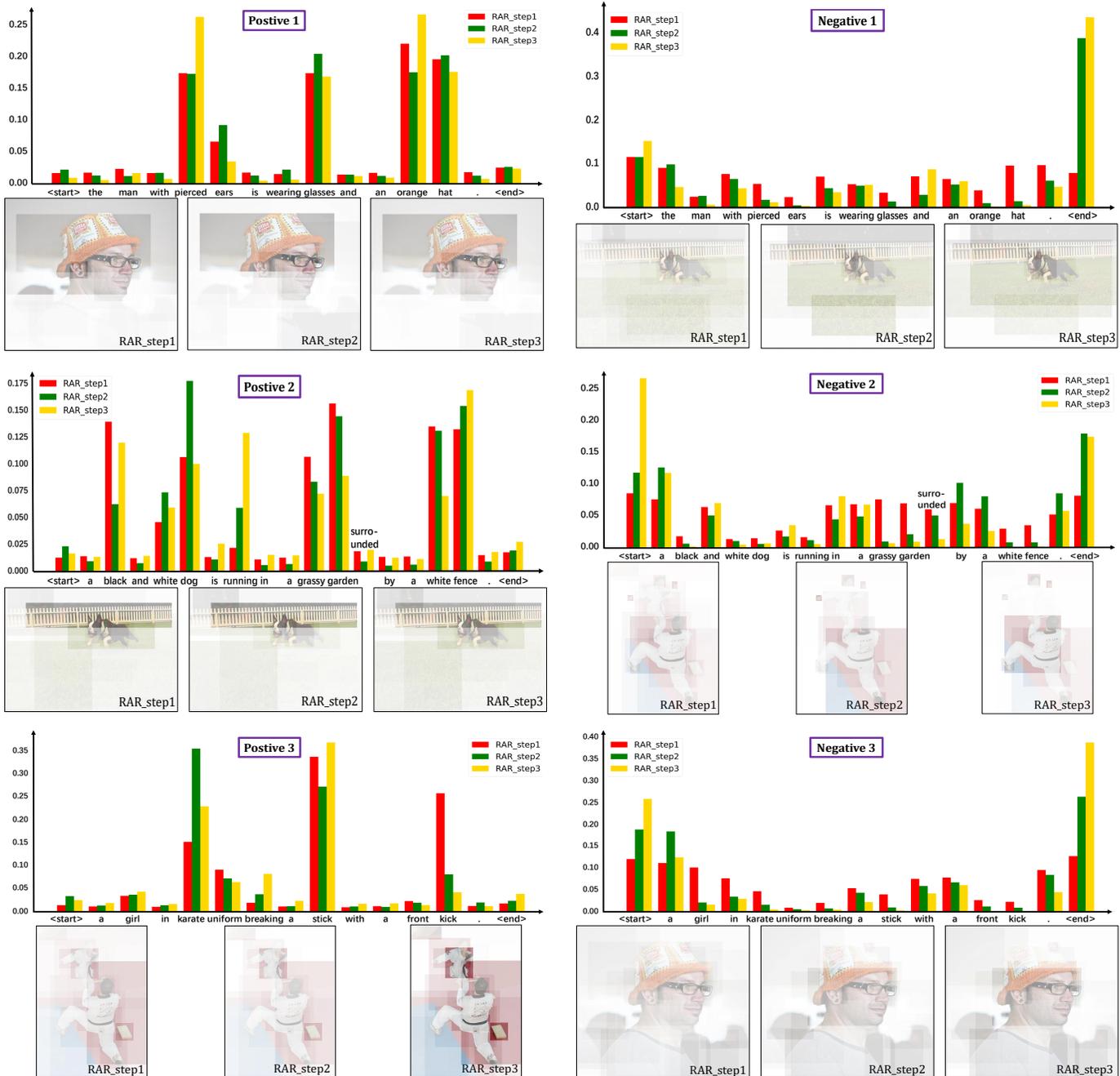


Fig. 8. Qualitative aggregation distribution by n-step RAR in the positive (right) and negative (left) pairs on Flickr30K. The histograms display the attention weights on word-based alignments with T2I attention while the images reflect the relative weights on region-based alignments with I2T attention.

obviously fail to measure feature channels and adapt itself to diverse words from different image-text pairs. Besides, even for the completely irrelevant image, the raw attention module still aligns the word with so-called "related regions" based on the cosine-like metrics and implicitly narrows the distances between the word and its related regions. In contrast, the RCR can make fine-grained adjustments with the prior alignments and refine the word-region correspondence progressively to produce larger gaps between matched and unmatched pairs.

Qualitative aggregation distribution of n-step RAR.

Fig. 8 illustrates the aggregation weights of word/region-attended alignments at the last step. We take Positive 1 as

an example of positive pairs. The RAR with T2I attention can highlight the discriminative alignments (*pierced ears, glasses, orange hat*) and abandon irrelevant ones (*the, with, is, etc.*), while with I2T attention, it can also capture salient regions mentioned in the text (*hat, ears, glasses*). Besides for negative pairs, the RAR with T2I attention tends to emphasize *<start>/<end>*-attended alignments which encode the overall image-text discrepancy. In terms of I2T attention, the aggregation distribution of image regions is relatively smooth to gain a more comprehensive prediction from the perspective of the image. As we can see, the RAR can selectively integrate important alignments and suppress less important ones.



Fig. 9. Qualitative T2I attention distribution and word-based cosine similarities with diverse semantics by n-step RCR in the positive (right) and negative (left) pairs on Flickr30K. The image and *_sim indicate the regions of interest and the corresponding cosine similarity according to the particular word.

Qualitative cross-attention distribution of n-step RCR.

Fig. 9 exhibits the regions of interest and corresponding similarities with respect to the words with various semantics. Note that the final image-text score by the RCR is also computed by averaging all the word-attended cosine similarities. Hence, the similarity between a word and integrated regions can reflect their correspondence quantitatively (A higher score means a higher correlation, and vice versa). We can observe that the RCR can refine the word-region interactions step by step and gradually draw the distance between diverse words and their related regions for positive pairs. Compared with *verbs* and *adjectives*, *nouns* are relatively easy to match for the Base model (T2I-SCAN [18]). When the RCR is introduced, the nouns/verbs/adjectives-based correspondences become more accurate and fine-grained. More importantly, semantics-based similarities in negative pairs are pulled down significantly, indicating that plug-in RCR learns the difference and relationship from the previous alignments and reweighs the channel-wise and word-wise attention factors to associate words with "completely irrelevant" regions in the latent space. By this means, the RCR can promote larger margins between positive and negative cross-modal pairs, and possess the greater capability to handle complex matching patterns.

VI. CONCLUSION AND FUTURE WORKS

In this paper, we proposed two regulators termed as Recurrent Correspondence Regulator (RCR) and Recurrent

Aggregation Regulator (RAR) to significantly facilitate the image-text matching process. Specifically, the RCR attempts to promote the cross-modal attention unit dynamically via learning more targeted attention factors, while the RAR aims to integrate the alignments progressively with plausible aggregation weights from holistic message feedback. The plug-and-play property enables them to seamlessly integrate into many existing approaches based on cross-modal interaction for achieving remarkable improvements, and more benefits can be obtained in a collaborative manner. Extensive experiments on MSCOCO and Flickr30K demonstrate the great superiority and broad applicability of our proposed approach. Beyond the above observations, we also attempt to apply our regulators to another branch [9], [13], [14], [17] focusing on single-modality representations without cross-modality interactions. Interestingly, 2-step RCR and 3-step RAR can improve the R@1 of SAEM [13] by 2.8/4.1% and 3.7/2.5% at two directions via gradually updating the last self-attention layer among regions and aggregating all instance features into a holistic feature respectively, reflecting the tremendous potential of our regulators. More efficient frameworks and application scenarios are one of our future research directions.

REFERENCES

- [1] S. Chen, Y. Zhao, Q. Jin, and Q. Wu, "Fine-grained video-text retrieval with hierarchical graph reasoning," in *CVPR*, 2020, pp. 10 635–10 644.
- [2] H. Wang, D. Xu, D. He, F. Li, Z. Ji, J. Han, and E. Ding, "Boosting video-text retrieval with explicit high-level semantics," in *ACMMM*, 2022, pp. 4887–4898.

- [3] S. Liu, H. Fan, S. Qian, Y. Chen, W. Ding, and Z. Wang, "Hit: Hierarchical transformer with momentum contrast for video-text retrieval," in *ICCV*, 2021, pp. 11 895–11 905.
- [4] H. Nam, J. Ha, and J. Kim, "Dual attention networks for multimodal reasoning and matching," in *CVPR*, 2017, pp. 2156–2164.
- [5] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, "Bottom-up and top-down attention for image captioning and visual question answering," in *CVPR*, 2018, pp. 6077–6086.
- [6] P. Wang, Q. Wu, J. Cao, C. Shen, L. Gao, and A. van den Hengel, "Neighbourhood watch: Referring expression comprehension via language-guided graph attention networks," in *CVPR*, 2019, pp. 1960–1968.
- [7] R. Zellers, Y. Bisk, A. Farhadi, and Y. Choi, "From recognition to cognition: Visual commonsense reasoning," in *CVPR*, 2019, pp. 6720–6731.
- [8] L. Wang, Y. Li, and S. Lazebnik, "Learning deep structure-preserving image-text embeddings," in *CVPR*, 2016, pp. 5005–5013.
- [9] F. Faghri, D. J. Fleet, J. R. Kiros, and S. Fidler, "VSE++: improving visual-semantic embeddings with hard negatives," in *BMVC*, 2018, p. 12.
- [10] Y. Huang, Q. Wu, C. Song, and L. Wang, "Learning semantic concepts and order for image and sentence matching," in *CVPR*, 2018, pp. 6163–6171.
- [11] J. Gu, J. Cai, S. R. Joty, L. Niu, and G. Wang, "Look, imagine and match: Improving textual-visual cross-modal retrieval with generative models," in *CVPR*, 2018, pp. 7181–7189.
- [12] Y. Song and M. Soleymani, "Polysemous visual-semantic embedding for cross-modal retrieval," in *CVPR*, 2019, pp. 1979–1988.
- [13] Y. Wu, S. Wang, G. Song, and Q. Huang, "Learning fragment self-attention embeddings for image-text matching," in *ACMMM*, 2019, pp. 2088–2096.
- [14] K. Li, Y. Zhang, K. Li, Y. Li, and Y. Fu, "Visual semantic reasoning for image-text matching," in *ICCV*, 2019, pp. 4653–4661.
- [15] S. Wang, R. Wang, Z. Yao, S. Shan, and X. Chen, "Cross-modal scene graph matching for relationship-aware image-text retrieval," in *WACV*, 2020, pp. 1497–1506.
- [16] H. Wang, Y. Zhang, Z. Ji, Y. Pang, and L. Ma, "Consensus-aware visual-semantic embedding for image-text matching," in *ECCV*, vol. 12369, 2020, pp. 18–34.
- [17] J. Chen, H. Hu, H. Wu, Y. Jiang, and C. Wang, "Learning the best pooling strategy for visual semantic embedding," in *CVPR*, 2021, pp. 15 789–15 798.
- [18] K. Lee, X. Chen, G. Hua, H. Hu, and X. He, "Stacked cross attention for image-text matching," in *ECCV*, vol. 11208, 2018, pp. 212–228.
- [19] C. Liu, Z. Mao, A. Liu, T. Zhang, B. Wang, and Y. Zhang, "Focus your attention: A bidirectional focal attention network for image-text matching," in *ACMMM*, 2019, pp. 3–11.
- [20] X. Wei, T. Zhang, Y. Li, Y. Zhang, and F. Wu, "Multi-modality cross attention network for image and sentence matching," in *CVPR*, 2020, pp. 10938–10947.
- [21] Z. Hu, Y. Luo, J. Lin, Y. Yan, and J. Chen, "Multi-level visual-semantic alignments with relation-wise dual attention network for image and text matching," in *IJCAI*, 2019, pp. 789–795.
- [22] Y. Wang, H. Yang, X. Qian, L. Ma, J. Lu, B. Li, and X. Fan, "Position focused attention network for image-text matching," in *IJCAI*, 2019, pp. 3792–3798.
- [23] T. Chen and J. Luo, "Expressing objects just like words: Recurrent visual embedding for image-text matching," in *AAAI*, 2020, pp. 10 583–10 590.
- [24] J. Wehrmann, C. Kolling, and R. C. Barros, "Adaptive cross-modal embeddings for image-text alignment," in *AAAI*, 2020, pp. 12 313–12 320.
- [25] Q. Zhang, Z. Lei, Z. Zhang, and S. Z. Li, "Context-aware attention network for image-text retrieval," in *CVPR*, 2020, pp. 3533–3542.
- [26] H. Chen, G. Ding, X. Liu, Z. Lin, J. Liu, and J. Han, "IMRAM: iterative matching with recurrent attention memory for cross-modal image-text retrieval," in *CVPR*, 2020, pp. 12 652–12 660.
- [27] C. Liu, Z. Mao, T. Zhang, H. Xie, B. Wang, and Y. Zhang, "Graph structured network for image-text matching," in *CVPR*, 2020, pp. 10 918–10 927.
- [28] H. Diao, Y. Zhang, L. Ma, and H. Lu, "Similarity reasoning and filtration for image-text matching," in *AAAI*, 2021, pp. 1218–1226.
- [29] Y. Huang, W. Wang, and L. Wang, "Instance-aware image and sentence matching with selective multimodal LSTM," in *CVPR*, 2017, pp. 7254–7262.
- [30] Z. Wang, X. Liu, H. Li, L. Sheng, J. Yan, X. Wang, and J. Shao, "CAMP: cross-modal adaptive message passing for text-image retrieval," in *ICCV*, 2019, pp. 5763–5772.
- [31] K. Lee, H. Palangi, X. Chen, H. Hu, and J. Gao, "Learning visual relation priors for image-text matching and image captioning with neural scene graph generators," *arXiv: 1909.09953*, 2019.
- [32] M. Karlsson and M. Covell, "Dynamic black-box performance model estimation for self-tuning regulators," in *ICAC*, 2005, pp. 172–182.
- [33] J. Xu, Y. Pan, X. Pan, S. C. H. Hoi, Z. Yi, and Z. Xu, "Regnet: Self-regulated network for image classification," *arXiv: 2101.00590*, 2021.
- [34] K. Thakkar, V. Paredes, and A. Hereid, "Adaptive feedback regulator for powered lower-limb exoskeleton under model uncertainty," *arXiv: 2104.11775*, 2021.
- [35] M. Abu-Khalaf, S. Karaman, and D. Rus, "Feedback from pixels: Output regulation via learning-based scene view synthesis," in *LADC*, vol. 144, 2021, pp. 828–841.
- [36] Y. Li, D. Zhang, and Y. Mu, "Visual-semantic matching by exploring high-order attention and distraction," in *CVPR*, 2020, pp. 12 783–12 792.
- [37] X. Dong, H. Zhang, L. Zhu, L. Nie, and L. Liu, "Hierarchical feature aggregation based on transformer for image-text matching," *TCSVT*, vol. 32, no. 9, pp. 6437–6447, 2022.
- [38] L. Qu, M. Liu, J. Wu, Z. Gao, and L. Nie, "Dynamic modality interaction modeling for image-text retrieval," in *SIGIR*, 2021, pp. 1104–1113.
- [39] K. Zhang, Z. Mao, Q. Wang, and Y. Zhang, "Negative-aware attention framework for image-text matching," in *CVPR*, 2022, pp. 15 661–15 670.
- [40] Y. Wang, T. Zhang, X. Zhang, Z. Cui, Y. Huang, P. Shen, S. Li, and J. Yang, "Wasserstein coupled graph learning for cross-modal retrieval," in *ICCV*, 2021, pp. 1813–1822.
- [41] J. Li, L. Niu, and L. Zhang, "Action-aware embedding enhancement for image-text retrieval," in *AAAI*, 2022, pp. 1323–1331.
- [42] H. Wang, D. He, W. Wu, B. Xia, M. Yang, F. Li, Y. Yu, Z. Ji, E. Ding, and J. Wang, "CODER: coupled diversity-sensitive momentum contrastive learning for image-text retrieval," in *ECCV*, 2022.
- [43] Z. Ji, K. Chen, and H. Wang, "Step-wise hierarchical alignment network for image-text matching," in *IJCAI*, 2021, pp. 765–771.
- [44] H. Zhang, Z. Mao, K. Zhang, and Y. Zhang, "Show your faith: Cross-modal confidence-aware network for image-text matching," in *AAAI*, 2022, pp. 3262–3270.
- [45] T. Tirer and R. Giryes, "Image restoration by iterative denoising and backward projections," *TIP*, vol. 28, no. 3, pp. 1220–1234, 2019.
- [46] Z. Zha, X. Yuan, J. T. Zhou, J. Zhou, B. Wen, and C. Zhu, "The power of triply complementary priors for image compressive sensing," in *ICIP*, 2020, pp. 983–987.
- [47] Z. Zha, B. Wen, X. Yuan, J. T. Zhou, J. Zhou, and C. Zhu, "Triply complementary priors for image restoration," *TIP*, vol. 30, pp. 5819–5834, 2021.
- [48] K. Zhang, Y. Li, W. Zuo, L. Zhang, L. V. Gool, and R. Timofte, "Plug-and-play image restoration with deep denoiser prior," *TPAMI*, vol. 44, no. 10, pp. 6360–6376, 2022.
- [49] Y. Chen, S. Wang, W. Zhang, and Q. Huang, "Less is more: Picking informative frames for video captioning," in *ECCV*, vol. 11217, 2018, pp. 367–384.
- [50] B. Wang, L. Ma, W. Zhang, and W. Liu, "Reconstruction network for video captioning," in *CVPR*, 2018, pp. 7622–7631.
- [51] L. Li, Z. Gan, Y. Cheng, and J. Liu, "Relation-aware graph attention network for visual question answering," in *ICCV*, 2019, pp. 10 312–10 321.
- [52] L. Chen, X. Yan, J. Xiao, H. Zhang, S. Pu, and Y. Zhuang, "Counterfactual samples synthesizing for robust visual question answering," in *CVPR*, 2020, pp. 10 797–10 806.
- [53] F. Shu, B. Chen, Y. Liao, S. Xiao, W. Sun, X. Li, Y. Zhu, J. Wang, and S. Liu, "Masked contrastive pre-training for efficient video-text retrieval," *arXiv: 2212.00986*, 2022.
- [54] R. Yan, M. Z. Shou, Y. Ge, A. J. Wang, X. Lin, G. Cai, and J. Tang, "Video-text pre-training with learned regions," *arXiv: 2112.01194*, 2021.
- [55] S. Ren, K. He, R. B. Girshick, and J. Sun, "Faster R-CNN: towards real-time object detection with region proposal networks," in *NIPS*, 2015, pp. 91–99.
- [56] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L. Li, D. A. Shamma, M. S. Bernstein, and L. Fei-Fei, "Visual genome: Connecting language and vision using crowdsourced dense image annotations," *IJCV*, vol. 123, no. 1, pp. 32–73, 2017.
- [57] M. Schuster and K. K. Paliwal, "Bidirectional recurrent neural networks," *TSP*, vol. 45, no. 11, pp. 2673–2681, 1997.
- [58] Z. Kuang, Y. Gao, G. Li, P. Luo, Y. Chen, L. Lin, and W. Zhang, "Fashion retrieval via graph reasoning networks on a similarity pyramid," in *ICCV*, 2019.
- [59] T. Lin, M. Maire, S. J. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: common objects in context," in *ECCV*, vol. 8693, 2014, pp. 740–755.

- [60] P. Young, A. Lai, M. Hodosh, and J. Hockenmaier, "From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions," *TACL*, vol. 2, pp. 67–78, 2014.
- [61] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR*, 2015.