

# Deep Boosting Learning: A Brand-new Cooperative Approach for Image-Text Matching

Haiwen Diao, Ying Zhang, Shang Gao, Xiang Ruan, Huchuan Lu

**Abstract**—Image-text matching remains a challenging task due to heterogeneous semantic diversity across modalities and insufficient distance separability within triplets. Different from previous approaches focusing on enhancing multi-modal representations or exploiting cross-modal correspondence for more accurate retrieval, in this paper we aim to leverage the knowledge transfer between peer branches in a boosting manner to seek a more powerful matching model. Specifically, we propose a brand-new Deep Boosting Learning (DBL) algorithm, where an anchor branch is first trained to provide insights into the data properties, with a target branch gaining more advanced knowledge to develop optimal features and distance metrics. Concretely, an anchor branch initially learns the absolute or relative distance between positive and negative pairs, providing a foundational understanding of the particular network and data distribution. Building upon this knowledge, a target branch is concurrently tasked with more adaptive margin constraints to further enlarge the relative distance between matched and unmatched samples. Extensive experiments validate that our DBL can achieve impressive and consistent improvements based on various recent state-of-the-art models in the image-text matching field, and outperform related popular cooperative strategies, e.g., Conventional Distillation, Mutual Learning, and Contrastive Learning. Beyond the above, we confirm that DBL can be seamlessly integrated into their training scenarios and achieve superior performance under the same computational costs, demonstrating the flexibility and broad applicability of our proposed method.

**Index Terms**—Image-text matching, Deep boosting learning, Deep cooperative learning, Deep metric learning.

## I. INTRODUCTION

With the explosion of multimedia volume in recent years, image-text matching [1], [2] has been a prevalent research topic, which efficiently bridges the gap between vision and language, and potentially benefits other multi-modal tasks such as video-text retrieval [3]–[5], referring expression [6], [7], and visual question answering [8], [9], etc. Despite years of efforts, image-text matching remains challenging because it entails not only recognizing hierarchical contents across modalities [10], [11], but also mapping diverse inputs into a comparable space to exploit semantic associations [12]–[14].

To explore a shared embedding space for cross-modal data, some works [15], [16] employ a hinge-based ranking loss function which forces each image/text to be closer to its positive text/image than all negatives within a mini-batch. Each triplet would be punished when the relative distance between

query-positive and query-negative pairs is less than a fixed margin. Though considering all pairs makes the optimization more stable, its sum-margin strategy treats all triplets equally during optimization, which would diminish the impact of valuable ones. To excavate more informative matching details, Faghri *et al.* [17] and Wei *et al.* [18] propose the max-margin and polynomial loss respectively to assign appropriate weights and highlight significant pairs from redundant pairs. However, the fixed distance margin for all triplets does not necessarily lead to good separability between the positive and negative samples. Hence, Zhao *et al.* [19] employs adaptive thresholds by computing the feature distances between text-to-text pairs as a reference, while Biten *et al.* [20] takes captioning metric (SPICE or CIDEr) as a measure and generates the semantic boundary via the language continuum of each caption. Besides, Zhou *et al.* [21] proposes the ladder loss with an inequality chain and adopts hierarchical margins for all triplets. However, they all consider an implicit and coarse relevance representation between each query and its candidates, resulting in an imprecise and inconsistent threshold constraint, and matching ambiguities between positive and negative pairs.

From the above perspective, we reconsider the key ingredients to fully exploit the potential of a matching network, i.e., constructing specialized guidance and seeking appropriate penalties. In this paper, we propose a novel Deep Boosting Learning (DBL) strategy, where an anchor branch is trained synchronously or asynchronously to provide explicit adaptive constraint for each triplet, in order to obtain a more powerful target branch. Different from previous metrics imposing explicit handcrafted penalties, the anchor branch learns the distance distribution and triplet relationship from data in advance, and would naturally gain an insight into the model properties and matching patterns, offering its twin target branch an explicit measurement to further enlarge distance separation and gain more discriminative feature metric. More specifically, the penalty of each triplet would be adjusted dynamically according to the similarity values predicted by the anchor branch, aiming to increase the association/separability between matched/unmatched pairs, as well as relax the tedious and costly configurations for robust constraint exploration. In this way, the target branch would capture a comprehensive picture of data and model characteristics, and translate the prior distance reference into more powerful matching capacities.

We notice that the proposed DBL is closely relevant to previous cooperative learning strategies including Conventional Distillation, Mutual Learning, and Contrastive Learning - they attempt to transfer prior knowledge and achieve better performance across multiple branches compared with

*Corresponding author: Huchuan Lu (lhchuan@dlut.edu.cn).* This work was supported by the National Natural Science Foundation of China under grant No. 62293540, 62293542. H. Diao, S. Gao are with Dalian University of Technology, China. (Email: diaohw@mail.dlut.edu.cn; gs940601k@gmail.com). Y. Zhang is with Tencent Company, China. (Email: yinggzhang@tencent.com). X. Ruan is with Tiwaki Company, Japan. (Email: ruanxiang@tiwaki.com).

independent learning. Specifically, they typically adopt a static pre-trained network [22]–[24], peer-teaching cohorts [25], or a slowly progressing encoder [26]–[28] as the anchor branch. To learn better probability prediction or feature representation, the target branch is encouraged to mimic the outputs of the anchor branch, including hard/soft pseudo-labels [25], [29], [30], absolute/relative relations [31]–[33], and feature similarities [27], [34], [35]. However, these approaches are originally designed for uni-modal tasks, and none of them take into account the heterogeneous semantic gap in cross-modal data. Besides, our DBL goes deeper into the margin knowledge in a peer-boosting manner beyond their peer-imitating ways to gain greater benefits under the same training and inference schemes, confirming the necessity of investigating effective cooperative strategies for cross-modal retrieval.

Our contributions are summarized as follows:

- We propose a novel Deep Boosting Learning (DBL) for peer-training strategy, which introduces an adaptive and explicit margin constraint for each triplet, and effectively generates the initiative distance separability between positive and negative pairs for image-text matching.
- Our DBL strategy can be widely applied to multiple training scenarios of related cooperative approaches, either as a post-processing step or in an online manner via collaborative or momentum synchronous updates.
- We validate the proposed DBL strategy with recent state-of-the-art works. Extensive experimental results on Flickr30K and MSCOCO datasets demonstrate the superiority and flexibility of our boosting strategy.

## II. RELATED WORK

### A. Image-Text Matching.

Image-text matching task targets retrieving images from the database with natural language queries, and vice versa. Research on this topic can be roughly divided into two aspects: **1) mono-modal representation.** To achieve this, some works [36]–[39] introduced graph reasoning networks to enhance the region and word features, while Wang *et al.* [40] utilized a constructed concept correlation to generate the consensus-aware embeddings. Besides, Li *et al.* [41], Chen *et al.* [42] and Zhu *et al.* [43] designed global memory bank, generalized pooling operator, and external space attention strategy respectively which effectively enhance the feature representation and facilitate mono-modal aggregation. Another set of methods [44]–[47] focuses on **2) cross-modal interaction.** For example, some approaches [48], [49] employed an iterative scheme with attention memory or regulator modules to recurrently refine region-word alignments, while several methods [39], [50]–[52] developed cross-modal correspondences to perform hierarchical matching with complex graph reasoning and high computational cost. Moreover, Zhang *et al.* [53] used the optimal boundary to explicitly and adaptively model the mismatched fragments and yield more accurate predictions. To evaluate the effectiveness and generalization of our proposed strategy, we apply our DBL to a series of representative works including conventional and pre-trained matching networks on the above two directions, and achieve

solid and consistent improvements on two benchmarks. Note that we also construct a concise but powerful baseline, which, though not our contribution, only serves as an insight into the mechanisms and comparisons of our DBL strategy.

### B. Deep Metric Learning.

Deep metric learning aims to map samples into a unified projection space, such that the similarities between positive pairs are higher than the ones between negative pairs. In the past few years, various loss functions for uni-modal retrieval tasks [54], [55] have been introduced, including triplet [54], quadruplet [55], lifted structure [56], N-pair [57], histogram [58], and Proxy-NCA [59], some of which have been extended for cross-modal matching. Wang *et al.* [15] proposed a two-way ranking loss by adapting the triplet loss for bi-directional retrievals, which has gained great popularity in multi-modal learning [60], [61]. Faghri *et al.* [17] introduced hard negative mining into the loss function, while Zhang *et al.* [62] developed the cross-modal projection loss, which minimizes the matching distribution between all pairs in a mini-batch. As mentioned before, the most related works are [20], [21], which regarded the relationship between text-to-text pairs as a reference and employed adaptive margin restrictions based on implicit semantic relevance degrees. Different from them, the DBL strategy automatically seeks an appropriate threshold according to the explicit distance within a triplet, and forces an adequate distance separation between matched and unmatched image-text pairs.

### C. Deep Cooperative Learning.

Deep cooperative learning is a typical peer-training strategy trading training efficiency for performance benefits, which brings extra training costs but no computational cost for inference. Specifically, conventional distillation trains a smaller student network to mimic the knowledge flow of a powerful yet static teacher, consisting of normalized probabilities [29], [30], network parameters [22], [23], [63], and feature relations [31], [32], while in mutual learning [25], the student cohorts imitate the predictions from each other, and their training processes are collaborative. Moreover, contrastive learning [26]–[28], [34] presents a promising way of unsupervised representation learning, and the key idea is constructing similar or dissimilar data examples and maximizing agreement between two encoder networks. In contrast, our DBL strategy starts with an anchor branch, first snooping on the prior relationships within triplets, and punishes a target branch via adaptive and adequate margin values. By this means, the latter can obtain a greater discriminative ability and achieve a better matching capability for single-branch learning. Besides, we experimentally validate that it can perform particularly well under multiple training scenarios of the above-mentioned approaches, reflecting the flexibility and wide suitability of our DBL strategy.

## III. METHODOLOGY

We first introduce a powerful single-branch network in detail, including feature extraction, cross-modal interaction, similarity prediction, and matching loss. We then elaborate on the detailed mathematics and training strategies of our DBL.

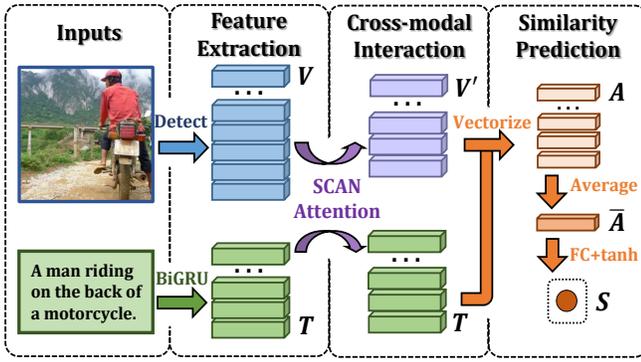


Fig. 1. Illustration of our single branch baseline. We adopt hard ranking loss [14] as a task-specific loss to supervise the training process.

### A. Single Branch Baseline

We build a simple and effective image-text matching baseline, which simply combines cross-attention module [44] and vectorized similarity representation [51], and is only used to analyze the mechanism and comparison of our DBL.

**Feature Extraction.** For each image, we first apply bottom-up attention [64] pretrained on Visual Genomes [65] to extract the top  $K$  region proposals with 2048-d features. Then, a fully-connected (FC) layer is utilized to map these features into 1024-d vectors  $\mathbf{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_K\} \in \mathbb{R}^{K \times 1024}$ . For a sentence with  $L$  words, we first encode them into 300-d word embeddings with random initialization, followed by a Bi-GRU [66] to integrate the bidirectional contextual information. Finally, we average the forward and backward hidden states at each time to get the word features  $\mathbf{T} = \{\mathbf{t}_1, \dots, \mathbf{t}_L\} \in \mathbb{R}^{L \times 1024}$ , and  $\mathbf{t}_l$  denotes the  $l$ -th word vector.

**Cross-modal Interaction.** We employ widely-used cross-modal attention [44] to capture region-word correspondence. Here, we take text-to-image attention as the backbone. To be specific, we first compute the cosine similarity matrix  $\mathbf{M} \in \mathbb{R}^{L \times K}$  for all region-word pairs, followed by a zero threshold and word-wise L2 normalization (norm). Then, we adopt the region-wise softmax function and integrate all the regions attended by each word as:

$$\begin{aligned} \mathbf{M} &= \text{norm}_{\mathbf{T}} \left( \left[ \mathbf{T} \mathbf{V}^\top \right]_+ \right), \\ \mathbf{V}' &= \text{softmax}_{\mathbf{V}} (\lambda \mathbf{M}) \mathbf{V}, \end{aligned} \quad (1)$$

where  $\lambda = 9$  following [44], and  $[x]_+ = \max(x, 0)$ . Note that  $\mathbf{V}' = \{\mathbf{v}'_1, \dots, \mathbf{v}'_L\} \in \mathbb{R}^{L \times 1024}$ , and  $\mathbf{v}'_l$  denotes the attended region features with respect to  $l$ -th word feature.

**Similarity Prediction.** As with [51], we first vectorize all the word-based alignments  $\mathbf{A} \in \mathbb{R}^{L \times 256}$  between  $\mathbf{T}$  and  $\mathbf{V}'$ , followed by the average operation to obtain one holistic alignment vector. Finally, we feed it into another FC layer and Tanh activation to output a scalar score:

$$\begin{aligned} \mathbf{A} &= \text{norm}(\mathbf{W}_1(|\mathbf{T} - \mathbf{V}'|^2) + \mathbf{b}_1), \\ \mathbf{S} &= \text{tanh}(\mathbf{W}_2(\bar{\mathbf{A}}) + \mathbf{b}_2), \end{aligned} \quad (2)$$

where  $|\cdot|^2$  denotes the element-wise square.  $\mathbf{W}_{\{\cdot\}}$  and  $\mathbf{b}_{\{\cdot\}}$  are learnable parameters, and  $\bar{\mathbf{A}} \in \mathbb{R}^{1 \times 256}$  represents word-wise average of  $\mathbf{A} \in \mathbb{R}^{L \times 256}$ , which indicates the similarity features attended by sentence words.

**Matching Loss.** Given a batch  $\mathcal{D} = \{(i_n, c_n)\}_{n=1}^N$  with  $N$  image-text pairs, the similarity outputs are denoted as  $\mathcal{S}_{\{\cdot, \cdot\}}$ . Note that  $\mathcal{S}_{i,c}$  and  $\mathcal{S}_{i,\hat{c}}/\mathcal{S}_{\hat{i},c}$  represent the matching scores of positive and negative pairs. Then, a hinge-based triplet ranking loss [16], [17] is widely used to guide optimization as:

$$\mathcal{L}_{raw} = \sum_{n=1}^N \ell(i_n, c_n). \quad (3)$$

**1) Sum-margin strategy.** It takes into account all possible combinations, ideally forcing all positive and negative samples to be separated by a margin value  $\gamma$ :

$$\mathcal{L}_{sum} = \sum_{\hat{c}} [\gamma + \mathcal{S}_{i,\hat{c}} - \mathcal{S}_{i,c}]_+ + \sum_{\hat{i}} [\gamma + \mathcal{S}_{\hat{i},c} - \mathcal{S}_{i,c}]_+, \quad (4)$$

where  $\hat{c}$  and  $\hat{i}$  are the negatives of  $i$  and  $c$ , and  $\gamma = 0.2$ .

**2) Max-margin strategy.** In contrast, the hard form only focuses on the nearest negatives ( $\hat{i}, \hat{c}$ ) in a mini-batch  $\mathcal{D}$ :

$$\mathcal{L}_{max} = [\gamma + \mathcal{S}_{i,\hat{c}} - \mathcal{S}_{i,c}]_+ + [\gamma + \mathcal{S}_{\hat{i},c} - \mathcal{S}_{i,c}]_+, \quad (5)$$

where  $\hat{c} = \arg \max_{d \neq c} \mathcal{S}_{i,d}$ ,  $\hat{i} = \arg \max_{j \neq i} \mathcal{S}_{j,c}$ .

**Discussion.** Although the latter can explore more informative details and effectively distinguish the confusing samples than the former, they both employ a handcrafted fixed threshold to restrict the relative distances between positive and negative pairs, resulting in an inadequate regularization that is easy for simple samples and hard for confusing ones. We empirically find that the network capability remains under-explored with such independent single-branch training. Hence, we propose the DBL strategy which carefully leverages peer knowledge to achieve greater matching capabilities.

### B. Deep Boosting Learning

The core idea is that, given the absolute or relative distance within triplets of the anchor branch, absolute or relative boosting strategies impose more compelling restraints on the corresponding distance for the target branch respectively, so that the target branch can obtain more suitable margin penalties and learn superior matching patterns not available in the anchor branch. To distinguish two branches, we denote the similarity scores of the target and anchor branch as  $\mathcal{S}_{\{\cdot, \cdot\}}^t$  and  $\mathcal{S}_{\{\cdot, \cdot\}}^a$ . The cumulative loss of boosting learning over training data  $\mathcal{D}$  is defined as:

$$\mathcal{L}_{boo} = \sum_{n=1}^N \ell'(i_n, c_n). \quad (6)$$

**Relative Boosting Strategy.** To fully make out the characteristics and relationships of each image-text pair, we first calculate relative distances between positive and negative pairs in the anchor branch, which serves as a prior and valuable insight into sample relationships of the original single branch. With the learned distance from the anchor branch, we introduce an adaptive margin for each triplet and impose more plausible restrictions when training the target branch, which we define as *Relative Sum (RS)* by:

$$\begin{aligned} \ell'_{RS} &= \sum_{\hat{c}} [\gamma + (\mathcal{S}_{i,c}^a - \mathcal{S}_{i,\hat{c}}^a) - (\mathcal{S}_{i,c}^t - \mathcal{S}_{i,\hat{c}}^t)]_+ \\ &\quad + \sum_{\hat{i}} [\gamma + (\mathcal{S}_{\hat{i},c}^a - \mathcal{S}_{i,c}^a) - (\mathcal{S}_{\hat{i},c}^t - \mathcal{S}_{i,c}^t)]_+. \end{aligned} \quad (7)$$

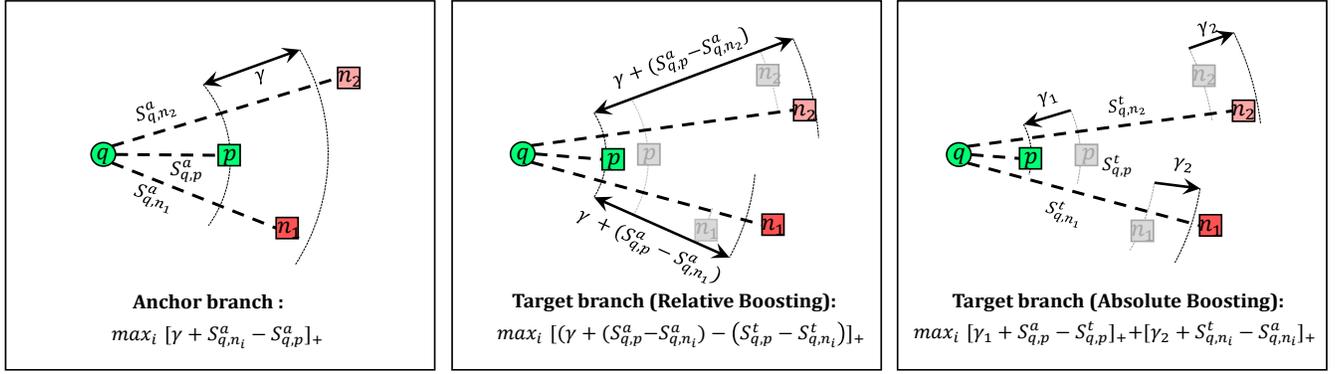


Fig. 2. Illustration of deep boosting learning. We first perform the anchor branch to obtain the absolute distance ( $S_{q,p}^a, S_{q,n_i}^a |_{i=1,2}$ ) between query and each candidate, and relative distance ( $S_{q,p}^a - S_{q,n_i}^a |_{i=1,2}$ ) within each triplet. Based on this prior knowledge, we assign the target branch with appropriate thresholds to further enlarge the variations between matched and unmatched image-text pairs.

With well-founded priori from the anchor branch, the goal of  $\ell'_{RS}$  is to implement adaptive penalties on the triplets, sufficiently pulling apart the relative distances while ensuring the stability of network convergence. However, a crucial caveat of the boosting strategy is the mining of hard negatives, as otherwise the training process will suffer from moderate negatives and will quickly stagnate. This is inspired by the analysis of minimizing a modified non-trivial loss function with uniform sampling in classification tasks [67], [68], identification tasks [54], [55], and multi-modal tasks [18], [69]. To emphasize hardest negatives ( $\check{c}, \check{i}$ ) for each positive pair ( $i, c$ ), we formulate *Relative Max (RM)* as:

$$\ell'_{RM} = [\gamma + (S_{i,c}^a - S_{i,\check{c}}^a) - (S_{i,c}^t - S_{i,\check{c}}^t)]_+ + [\gamma + (S_{i,c}^a - S_{i,c}^a) - (S_{i,c}^t - S_{i,c}^t)]_+, \quad (8)$$

where  $\check{i} = \arg \max_{\check{i}} (S_{i,c}^t - S_{i,\check{c}}^t)$  and  $\check{c} = \arg \max_{\check{c}} (S_{i,\check{c}}^t - S_{i,\check{c}}^a)$ . It is worth noting that ( $\check{i}, \check{c}$ ) represent the unmatched samples where the relative distances in the target branch are the toughest to push away based on the ones in the anchor branch, rather than the most confusing negatives ( $\hat{i}, \hat{c}$ ) in the target branch itself. At this point, the training difficulty of the target branch is relatively higher aiming to further capture discriminative matching details and improve the quality of the learned metrics. Even with its effectiveness, we argue that the relative boosting strategy does not specify how close the positive pairs are and how far the negative pairs are, inevitably rendering an insufficient exploration of positive-pair intimacy and negative-pair alienation.

**Absolute Boosting Strategy.** Based on the above observation, we explicitly normalize the absolute distances by comparing positive or negative pairs respectively between target and anchor branches. Hence, in contrast to the relative strategy, we directly adjust two new adaptive and explicit margins for the matched and unmatched pairs to pull the former closer, and meanwhile push the latter farther from each other in the target branch. Translating this statement into equation, we define *Absolute Sum (AS)* as:

$$\ell'_{AS} = \sum_{\check{c}} ([\gamma_1 + S_{i,c}^a - S_{i,c}^t]_+ + [\gamma_2 + S_{i,\check{c}}^t - S_{i,\check{c}}^a]_+) + \sum_{\check{i}} ([\gamma_1 + S_{i,c}^a - S_{i,c}^t]_+ + [\gamma_2 + S_{i,\check{i}}^t - S_{i,\check{i}}^a]_+), \quad (9)$$

where  $\gamma_1 = \alpha\gamma$ ,  $\gamma_2 = \gamma - \alpha\gamma$  (consistent with Eq. (4)(5)(7)). As described above, we also exploit hard negative mining to discover the hidden details between image regions and text words, and produce larger gaps between positives and negatives. The *Absolute Max (AM)* can be formulated as:

$$\ell'_{AM} = [\gamma_1 + S_{i,c}^a - S_{i,\check{c}}^t]_+ + [\gamma_2 + S_{i,\check{c}}^t - S_{i,\check{c}}^a]_+ + [\gamma_1 + S_{i,c}^a - S_{i,c}^t]_+ + [\gamma_2 + S_{i,c}^t - S_{i,c}^a]_+, \quad (10)$$

where the hardest negatives ( $\check{i}, \check{c}$ ) are mathematically equivalent to the ones of Eq. (8) in the same mini-batch. Different from the relative loss function that only requires the distance difference between anchor and target branches to be less than a unified margin, the absolute loss function attempts to simultaneously impose explicit and tighter penalties on absolute distances of positive and negative pairs. By this means, the latter can further enhance the discriminative power of the target branch, and develop the optimal feature and distance metric jointly for image-text matching.

**Discussion.** We utilize the relative strategy to only supervise the relative distance within triplets, while the absolute strategy further constrains the absolute distance among each pair.

**1) Relative vs. Absolute.** To take  $\ell'_{RM}$  and  $\ell'_{AM}$  as an example, we derive their connections by the formulas:

$$\begin{aligned} & [\gamma + (S_{i,c}^a - S_{i,\check{c}}^a) - (S_{i,c}^t - S_{i,\check{c}}^t)]_+ \\ &= [(\gamma_1 + S_{i,c}^a - S_{i,c}^t) + (\gamma_2 + S_{i,\check{c}}^t - S_{i,\check{c}}^a)]_+ \quad (11) \\ &\leq [\gamma_1 + S_{i,c}^a - S_{i,c}^t]_+ + [\gamma_2 + S_{i,\check{c}}^t - S_{i,\check{c}}^a]_+, \end{aligned}$$

$$\begin{aligned} & [\gamma + (S_{i,c}^a - S_{i,c}^a) - (S_{i,c}^t - S_{i,c}^t)]_+ \\ &= (\gamma_1 + S_{i,c}^a - S_{i,c}^t) + (\gamma_2 + S_{i,c}^t - S_{i,c}^a)]_+ \quad (12) \\ &\leq [\gamma_1 + S_{i,c}^a - S_{i,c}^t]_+ + [\gamma_2 + S_{i,c}^t - S_{i,c}^a]_+. \end{aligned}$$

Combining Eq. (11) and (12), we can obtain the inequality relations between two boosting strategies as follows:

$$\ell'_{RM} \leq \ell'_{AM}, \quad (13)$$

where the equality holds if and only if three items ( $\gamma_1 + S_{i,c}^a - S_{i,c}^t$ ), ( $\gamma_2 + S_{i,\check{c}}^t - S_{i,\check{c}}^a$ ), ( $\gamma_2 + S_{i,c}^t - S_{i,c}^a$ ) share the same signs. Similarly,  $\ell'_{RS} \leq \ell'_{AS}$ , confirming that the absolute form can impose tighter constraints, and produce more compact distances than the relative one (See Sec. IV-C).

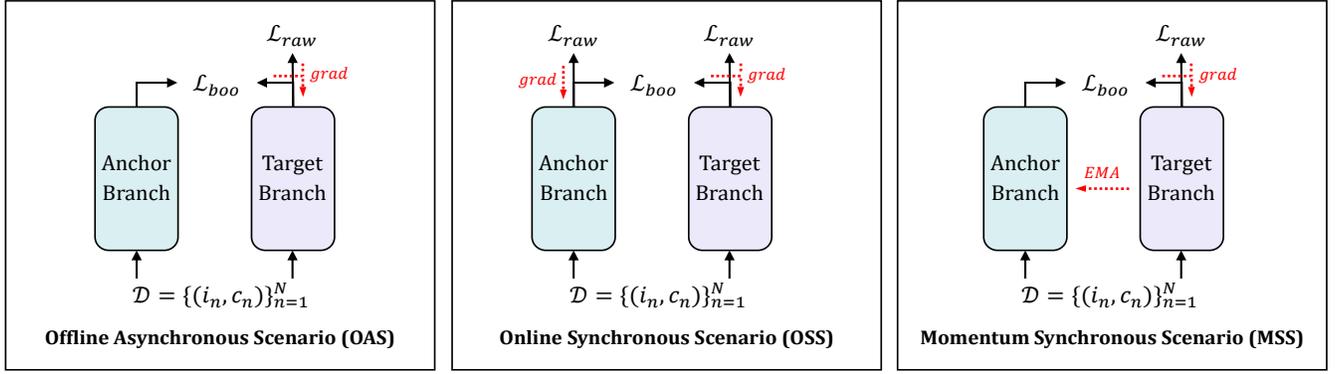


Fig. 3. Illustration of multiple training scenarios. OAS adopts two-stage training scheme as with Conventional Distillation, while OSS and MSS employ one-stage parallel training scenario as with Mutual Learning and Contrastive Learning, respectively. Notably, we only verify the target branch on the validation set, and utilize the model with the best RSUM to perform prediction on the test set.

**2) Fixed  $\gamma$  vs. Soft  $\gamma^{SA}$ .** The standard configurations in Eq. (7)(8)(9)(10) recommend RM and AM with fixed  $\gamma$ , meaning that the distance metric learned by the target branch has a consistent  $\gamma$  penalty based on the ones learned by the anchor branch. On the other hand, there is a value range of all available measures used for boosting strategies, e.g. the absolute distance of positive/negative pairs  $\mathcal{S}_{i,c}^a/\mathcal{S}_{i,\tilde{c}}^a/\mathcal{S}_{i,c}^a \in [-d_y, d_y]_{d_y=1}$  in Eq. (10), and the relative distance within each triplet  $(\mathcal{S}_{i,c}^a - \mathcal{S}_{i,\tilde{c}}^a)/(\mathcal{S}_{i,c}^a - \mathcal{S}_{i,c}^a) \in [-d_x, d_x]_{d_x=2}$  in Eq. (8). In other words, each type of the above distances theoretically has a boosting extreme namely **Theoretical Maximum (TM)**, some of which are less than the predefined  $\gamma$ . Combining these two findings, we propose a soft  $\gamma^{SA}$  formulation namely **Soft Adaptation (SA)** which varies exponentially corresponding to predicted distance from the anchor branch.

For  $\gamma$  in RM, we initialize  $\gamma^{SA}(x)$  that corresponds to the output score from the anchor branch as follows:

$$\gamma^{SA}(x) = \frac{2\gamma}{1 + e^{\epsilon(x-d_x)}} - \gamma, \quad (14)$$

where  $\epsilon$  is a smooth value that controls the sharpness of the curve  $\gamma^{SA}(x)$ .  $x$  indicates the inferred relative distance within each triplet from anchor branch. To maximize margin penalty without exceeding the extreme, we obtain an equation as:

$$\left. \frac{\partial \gamma^{SA}(x)}{\partial x} \right|_{x=d_x} = \left. \frac{\partial \gamma^{TM}(x)}{\partial x} \right|_{x=d_x}. \quad (15)$$

Given the above equation, we then compute the  $\epsilon$  value. The derivation is as follows:

$$\begin{aligned} \left. \frac{\partial \gamma^{SA}(x)}{\partial x} \right|_{x=d_x} &= \left. \frac{-2\gamma\epsilon e^{\epsilon(x-d_x)}}{(1 + e^{\epsilon(x-d_x)})^2} \right|_{x=d_x} = \frac{-\gamma\epsilon}{2}, \\ \left. \frac{\partial \gamma^{TM}(x)}{\partial x} \right|_{x=d_x} &= \left. \frac{\partial (d_x - x)}{\partial x} \right|_{x=d_x} = -1. \end{aligned} \quad (16)$$

Hence,  $\epsilon = \frac{2}{\gamma}$ . The  $\gamma^{SA}(x)$  for RM can be formulated as:

$$\gamma^{SA}(x) = \frac{2\gamma}{1 + e^{\frac{2}{\gamma}(x-d_x)}} - \gamma = \frac{\gamma - \gamma e^{\frac{2}{\gamma}(x-d_x)}}{1 + e^{\frac{2}{\gamma}(x-d_x)}}, \quad (17)$$

and we redefine the *Relative Max (RM)* as:

$$\begin{aligned} \ell'_{RM} &= [\gamma^{SA}(\mathcal{S}_{i,c}^a - \mathcal{S}_{i,\tilde{c}}^a) + (\mathcal{S}_{i,c}^a - \mathcal{S}_{i,\tilde{c}}^a) - (\mathcal{S}_{i,c}^t - \mathcal{S}_{i,\tilde{c}}^t)]_+ \\ &\quad + [\gamma^{SA}(\mathcal{S}_{i,c}^a - \mathcal{S}_{i,c}^a) + (\mathcal{S}_{i,c}^a - \mathcal{S}_{i,c}^a) - (\mathcal{S}_{i,c}^t - \mathcal{S}_{i,c}^t)]_+. \end{aligned} \quad (18)$$

Similarly for  $\gamma_1, \gamma_2$  in AM, we obtain  $\gamma_1^{SA}(y), \gamma_2^{SA}(y)$  as:

$$\begin{aligned} \gamma_1^{SA}(y) &= \frac{2\gamma_1}{1 + e^{\frac{2}{\gamma_1}(y-d_y)}} - \gamma_1 = \frac{\gamma_1 - \gamma_1 e^{\frac{2}{\gamma_1}(y-d_y)}}{1 + e^{\frac{2}{\gamma_1}(y-d_y)}}, \\ \gamma_2^{SA}(y) &= \frac{2\gamma_2}{1 + e^{\frac{2}{\gamma_2}(y+d_y)}} - \gamma_2 = \frac{\gamma_2 - \gamma_2 e^{\frac{2}{\gamma_2}(y+d_y)}}{1 + e^{\frac{2}{\gamma_2}(y+d_y)}}, \end{aligned} \quad (19)$$

and we reformulate the *Absolute Max (AM)* as:

$$\begin{aligned} \ell'_{AM} &= [\gamma_1^{SA}(\mathcal{S}_{i,c}^a) + \mathcal{S}_{i,c}^a - \mathcal{S}_{i,c}^t]_+ + [\gamma_2^{SA}(\mathcal{S}_{i,\tilde{c}}^a) + \mathcal{S}_{i,\tilde{c}}^a - \mathcal{S}_{i,\tilde{c}}^t]_+ \\ &\quad + [\gamma_1^{SA}(\mathcal{S}_{i,c}^a) + \mathcal{S}_{i,c}^a - \mathcal{S}_{i,c}^t]_+ + [\gamma_2^{SA}(\mathcal{S}_{i,c}^a) + \mathcal{S}_{i,c}^a - \mathcal{S}_{i,c}^t]_+. \end{aligned} \quad (20)$$

Experiments in Sec. IV-C show that the soft one displays stronger abilities of image retrieval but slightly attractive promotions on sentence retrieval, while the fixed one obtains the optimal balance between bidirectional retrievals. In summary, we recommend the fixed one as the vanilla boosting strategy.

### C. Multiple Training Scenarios

The collaboration of target and anchor branches is flexible. Following Conventional Distillation [29], [31], [32], [70], Mutual Learning [25], and Contrastive Learning [26]–[28], we adopt three popular and corresponding scenarios namely Offline Asynchronous, Online Synchronous and Momentum Synchronous Scenarios to validate the effectiveness and generalization of our boosting strategy.

**Conventional Distillation.** Traditional knowledge distillation [29]–[32] is a mechanism where the target (student) branch learns to match the results of the static anchor (teacher) branch, parameterized by  $\theta_t$  and  $\theta_a$  respectively:

$$\min_{\theta_t} \mathcal{L}_{kd}(F_{\theta_t}(x), F_{\theta_a}(x)), \quad (21)$$

where  $\mathcal{L}_{kd}$  consists of logit/distance/angle-wise distillation loss functions which help the target branch transfer the powerful knowledge from the pre-trained anchor branch. To

evaluate it, we first obtain a strong anchor branch by task-specific loss function  $\mathcal{L}_{raw}$ . As a post-processing step, the training procedure of the target branch is then supervised by  $\mathcal{L}_{raw} + \mathcal{L}_{boo}$  which penalizes the proximities between two branches and agitates the latter to gain the better matching ability, denoted as *Offline Asynchronous Scenario (OAS)*.

**Mutual Learning.** Unlike general knowledge distillation, the branch cohorts are updated jointly and collaboratively in deep mutual learning which does not rely on prior knowledge. For example, DML [25] attempts to optimize the cohorts by bringing their probability estimates closer and minimizing their discrepancies during the learning progress. However, our proposed approach aims to implicitly push aside the peers' distributions and exploit more powerful paradigms beyond mimicry. Hence, simultaneously using boosting strategy for each branch may lead to unclear training objectives and unstable optimizing paths, and ultimately converge to locally mediocre solutions for peer training. Therefore, we randomly initialize two branches where only the target branch is updated under the guidance of the boosting strategy by the anchor branch, denoted as *Online Synchronous Scenario (OSS)*.

**Contrastive Learning.** To avoid model collapse, several works focus on contrastive loss [71], inconsistent structure [27], [35], clustering constraint [72], [73], and momentum encoder [26], [28]. Likewise, they do not require a pre-trained network given a priori. Inspired by the momentum encoder [26], we update the anchor parameters  $\theta_a$  with an exponential moving average (EMA) of the target parameters  $\theta_t$ . Formally, we update  $\theta_a$  by:

$$\theta_a \leftarrow \beta\theta_a + (1 - \beta)\theta_t, \quad (22)$$

where  $\beta$  follows a cosine schedule [27], [28] from 0.99995 to 1 during training process. Only the target branch is updated by back-propagation of  $\mathcal{L}_{raw} + \mathcal{L}_{boo}$ , and the dynamic anchor branch progressively provides a consistent reference of higher quality and hence, we have no need of maintaining the queue dictionary of data samples or introducing data augmentation to form positive pairs. Note that this cooperative strategy serves as a standard operation similar to Polyak-Ruppert averaging with exponential decay [74], [75], which is denoted as *Momentum Synchronous Scenario (MSS)*.

**Discussion.** Deep cooperative learning is a popular technique trading extra training consumption for performance gains, and only utilizing the target branch for prediction under the above training scenarios. OAS adopts a two-stage training scheme as with DR [30] and RKD [32], while OSS and MSS train two branches simultaneously at one stage as with DML [25] and DINO [28] respectively. Compared with OSS, OAS and MSS require no gradients for the anchor branch during the cooperative process. We directly employ  $\mathcal{L}_{raw}$  and  $\mathcal{L}_{boo}$  with a 1:1 contribution to train the target branch, which has achieved steady and consistent improvements in all experiments without complex manual tuning. For a fair comparison, we ensure the same training and inference expenses, and validate that MSS can obtain great benefits with both slight training time and memory costs in TABLE II.

## IV. EXPERIMENTS

We first introduce the detailed training settings. Then, we report the cooperation and comparison with recent works and some popular strategies. After that, we investigate the configurations and analyses of our proposed DBL. Finally, we visualize some illustrations of bidirectional retrieval examples.

### A. Datasets and Settings

**Benchmark Datasets.** We evaluate our proposed approaches on two benchmark datasets: Flickr30K [76] and MSCOCO [77]. Each image of these two datasets is annotated with five corresponding captions. For Flickr30K, we adopt the standard split [1] and divide the dataset into 29,000 training images, 1,000 validation images, and 1,000 testing images. For MSCOCO, we follow [17], [44] to utilize 113,287 images for training, 5,000 images for validation and 5,000 images for testing. The 1k evaluation result is computed by averaging over 5 folds of 1K test images on MSCOCO.

**Evaluation Metrics.** We adopt Recall@ $\kappa$  ( $R@k$ ) and sum (RSUM) of  $R@1$ ,  $R@5$ , and  $R@10$  in two directions for evaluation, where  $R@k$  indicates the percentage of queries whose correct response is included in the top- $k$  candidates. Since  $R@k$  evaluation only cares about the first groundtruth retrieved in the top- $k$  results, we introduce the mean distance (MD) between positive and negative candidates for a better illustration of the similarity separability.

**Implementation Details.** Inspired by similarity representations [50]–[52], our baseline is an improved version of SCAN [44] with the vectorized similarity [51] instead of cosine distance. Following them, we extract  $K=36$  salient regions by bottom-up attention [64] for each image, and map 300-d word embeddings with random initialization into 1024-d features by Bi-GRU. We set the dimension of alignment vectors to be 256 with the inversed temperature  $\lambda$  as 9. The margin value  $\gamma$  and the proportion  $\alpha$  are set as 0.2 and 0.5 respectively. We train our method with 20 and 40 epochs on MSCOCO and Flickr30k dataset respectively during all the training scenarios, and set the initial learning rate as 0.0002 for 10 and 30 epochs, and decay it by 0.1 for the rest epochs.

### B. Quantitative Results

**Cooperation with Multiple Models.** TABLE I lists the applications on two types of image-text matching architectures, including Embedding-based (VSRN [36], ESA [43], CLIP [78]) and Interaction-based (BFAN [45], SGRAF [51], NAAF [53], OSCAR [79]) methods. Considering the diverse settings of these methods, it is essential to mention that we uniformly adopt the setups of BiGRU and a single model based on VSRN, ESA, BFAN, and NAAF. In particular, we utilize BFAN with equal attention and similarity representation, CLIP with ViT-L/14@336px encoder, and OSCAR with base BERT as references. Limited by current resources, we apply DBL under OSS for relatively smaller VSRN, ESA, BFAN, SGRAF, and NAAF, and under OAS for larger pre-trained CLIP and OSCAR. We keep their original loss functions, training settings, and model configurations. It is worth noting that we only utilize the target branch for performance validation.

TABLE I

COOPERATION WITH MULTIPLE REPRESENTATIVE MODELS INCLUDING EMBEDDING-BASED AND INTERACTION-BASED LEARNING ON FLICKR30K AND MSCOCO. WE RE-IMPLEMENT THESE METHODS WITH THEIR PUBLICLY AVAILABLE CODE. THE BEST RESULTS OF THE RSUM ARE MARKED IN **BOLD**.

Method	Flickr30K					MSCOCO 1K					MSCOCO 5K					
	Sentence Retrieval		Image Retrieval		RSUM	Sentence Retrieval		Image Retrieval		RSUM	Sentence Retrieval		Image Retrieval		RSUM	
	R@1	R@5	R@1	R@5		R@1	R@5	R@1	R@5		R@1	R@5	R@1	R@5		
Embedding Learning	<b>VSRN [36]</b>	70.2	89.4	53.2	78.0	471.2	74.2	94.1	60.6	88.3	509.2	50.3	79.4	37.6	68.5	403.3
	+RM (OSS)	72.1	90.3	54.8	78.6	476.3	75.3	94.6	61.5	89.0	512.8	51.6	79.9	38.8	69.5	407.4
	+AM (OSS)	72.8	90.4	55.0	78.9	<b>477.1</b>	75.2	94.8	61.8	89.2	<b>513.2</b>	51.7	80.0	39.4	69.7	<b>408.6</b>
	<b>ESA [43]</b>	82.3	95.8	61.2	86.0	514.5	79.2	96.4	63.5	90.8	524.9	58.0	84.8	41.2	71.3	429.3
	+RM (OSS)	83.4	96.1	61.9	86.4	517.3	80.0	96.5	63.7	91.1	526.3	58.6	85.0	41.5	72.0	431.2
	+AM (OSS)	83.2	96.2	62.2	86.5	<b>517.5</b>	80.1	96.5	63.8	91.2	<b>526.7</b>	58.8	85.2	41.6	72.0	<b>431.8</b>
Embedding Learning	<b>CLIP [78]</b>	92.0	99.4	77.8	95.0	561.1	83.2	96.4	68.3	91.3	534.8	67.5	88.3	49.4	75.1	457.8
	+RM (OAS)	93.0	99.5	79.1	95.3	<b>563.8</b>	83.9	96.8	68.7	91.6	536.5	68.2	88.6	49.8	75.3	459.7
	+AM (OAS)	92.8	99.5	78.9	95.2	563.3	84.1	97.0	68.8	91.6	<b>536.7</b>	68.3	88.8	50.1	75.4	<b>460.1</b>
	<b>BFAN [45]</b>	70.3	91.6	53.2	78.5	474.2	75.5	94.3	60.8	88.0	510.6	54.6	82.1	38.8	68.7	413.8
	+RM (OSS)	72.8	93.3	55.5	79.3	<b>483.5</b>	77.9	95.8	62.4	89.2	<b>518.7</b>	56.4	83.6	40.7	69.8	<b>421.9</b>
	+AM (OSS)	73.9	92.6	55.3	78.6	481.0	77.7	95.7	62.4	89.1	518.0	55.7	83.6	40.4	70.0	420.9
Interaction Learning	<b>SGRAF [51]</b>	78.1	94.4	58.2	83.1	500.2	79.3	96.2	63.5	90.3	523.8	58.1	85.0	41.8	71.3	429.6
	+RM (OSS)	80.1	94.7	59.6	83.5	504.9	79.8	96.7	64.0	90.5	525.5	59.2	84.9	42.3	71.6	432.2
	+AM (OSS)	79.8	95.4	60.7	84.0	<b>507.2</b>	80.5	96.6	64.3	90.6	<b>526.5</b>	59.8	85.2	42.5	71.5	<b>433.1</b>
	<b>NAAF [53]</b>	78.3	96.1	59.6	84.4	506.6	77.8	95.8	62.5	89.6	519.5	56.3	84.1	40.9	70.1	422.9
	+RM (OSS)	78.6	96.2	60.1	84.4	507.8	78.4	96.1	63.1	89.6	521.0	56.8	84.2	41.3	70.2	<b>424.3</b>
	+AM (OSS)	79.1	96.4	60.3	84.6	<b>509.1</b>	78.8	95.9	62.9	89.5	<b>521.4</b>	57.2	84.0	41.1	69.8	424.0
Interaction Learning	<b>OSCAR [79]</b>	-	-	-	-	-	88.4	99.1	75.7	95.2	556.5	70.0	91.1	54.0	80.8	479.9
	+RM (OAS)	-	-	-	-	-	88.8	98.9	75.8	95.4	557.1	71.0	91.0	54.5	80.9	481.2
	+AM (OAS)	-	-	-	-	-	88.8	99.0	76.1	95.5	<b>557.4</b>	70.9	91.1	54.8	81.0	<b>481.4</b>

*Visual Semantic Reasoning (VSRN)* [36] performs graph reasoning to generate visual features with semantic relationships.

*External Space Attention Aggregation (ESA)* [43] enables element-wise attention for discriminative information and adaptive feature aggregation at the dimensional level.

*Transferable Visual Models From Natural Language Supervision (CLIP)* [78] learns dual powerful image and text encoders on massive image-text pairs collected from the internet.

*Bidirectional Focal Attention (BFAN)* [45] enhances the region-word correspondence via a focal attention unit, and integrates all region-based and word-based similarities.

*Similarity Graph Reasoning and Attention Filtration (SGRAF)* [51] designs word-based similarity function, followed by graph reasoning and attention filtration modules.

*Negative-aware Attention Framework (NAAF)* [53] learns the optimal boundary to explicitly model matched and mismatched fragments, which yield the final score together.

*Object-Semantics Aligned Pre-training (OSCAR)* [79] utilizes object tags detected in images as a bridge to considerably alleviate cross-modality gap and alignment burden.

As can be clearly seen, they all benefit from the DBL strategy at all evaluation metrics for both relative and absolute manners. For R@1 at sentence/image retrieval on Flickr30K, it gains a maximal boost of 2.6/1.8% (VSRN), 1.1/1.0% (ESA), 3.6/2.3% (BFAN), 2.0/2.5% (SGRAF), and 0.8/0.7% (NAAF) respectively. The impressive and consistent improvements are also shown on MSCOCO 1K and 5K test sets, which well display strong capability and broad applicability regardless of network frameworks. Besides for pre-trained CLIP, our DBL improves R@1 by at most 1.0/1.3% and 0.8/0.7% on Flickr30K and MSCOCO5K, while based on OSCAR, our DBL still obtains a maximum R@1 boost of 1.0/0.8% on challenging MSCOCO5K. Notably, we have also attempted to apply DBL to GPO [42] and found that its warm-up process causes drastic changes in the loss magnitudes (1300-400, 48-

TABLE II

COMPARISON WITH COOPERATIVE LEARNING ON FLICKR30K. \* INDICATES THE INVERSE KL DIVERGENCE OF DISTILLATION LOSS.

Method	Sentence Retrieval				Image Retrieval		Scenarios		
	R@1	R@5	R@1	R@5	MD	Memory	Time		
Baseline	75.8	93.4	56.5	81.4	0.91	100%	100%		
DR [30]	76.8	93.7	57.4	81.8	1.25			OAS	
RKD [32]	77.6	93.8	56.8	81.9	0.74				
<b>RM (OAS)</b>	79.1	<b>94.6</b>	58.5	<b>83.3</b>	1.37			+12%	+120%
<b>AM (OAS)</b>	<b>79.7</b>	<b>94.6</b>	<b>58.7</b>	<b>83.2</b>	<b>1.58</b>				
DML [25]	77.7	93.6	57.4	82.8	0.94			OSS	
<b>RM (OSS)</b>	<b>79.3</b>	<b>95.4</b>	<b>59.1</b>	83.0	1.21				
<b>AM (OSS)</b>	79.0	94.9	58.5	<b>83.1</b>	<b>1.56</b>			+100%	+73%
DINO [28]	76.5	93.4	58.5	83.2	0.89			MSS	
DINO* [28]	77.9	93.9	58.8	83.5	1.15				
<b>RM (MSS)</b>	78.4	94.4	59.1	<b>83.6</b>	1.45			+11%	+18%
<b>AM (MSS)</b>	<b>79.4</b>	<b>94.7</b>	<b>59.7</b>	83.5	<b>1.81</b>				

14), making it difficult to balance the weights between GPO and DBL. Without a warm-up strategy, GPO with RM and AM can improve the bidirectional R@1, RSUM from 75.3/56.0%, 493.2% to 76.8/57.3%, 499.6% and 77.2/56.7%, 500.3% on Flickr30K, verifying consistent boosts and robustness of our DBL strategy. In conclusion, it is inadequate for these works supervised by [17] to distinguish the relations within triplets, and our DBL can assign adaptive and targeted margins and produce optimal matching patterns between image and text.

**Comparison with Deep Cooperative Learning.** For fairness, we compare DBL with some typical strategies under their original training settings, including DR (softrank), RKD (distance), DML, and DINO (ema). We extend them with some refinements to be more suitable for cross-modal tasks. For the first two works, we directly utilize the predicted similarity scores for knowledge distillation, while for the last two methods, there are no category labels to supply the imitation process of probability outputs. Hence, we introduce the cross-modal projection [62] that treats each image-text pair

TABLE III

COMPARISON WITH METRIC LEARNING ON FLICKR30K. HERE WE REPORT MSS DUE TO ITS CLOSE CONSUMPTION WITH SINGLE BRANCH.

Loss	Sentence Retrieval		Image Retrieval		MD
	R@1	R@5	R@1	R@5	
VSE0 [15]	74.4	92.8	54.3	81.2	0.58
VSE++ [17]	75.8	93.4	56.5	81.4	0.91
SAM [20]	76.1	93.2	57.0	82.9	0.88
CMPM [62]	76.8	95.1	57.8	82.8	1.12
MPL [18]	77.1	94.3	57.4	82.3	0.96
<b>RM (MSS)</b>	78.4	94.4	59.1	<b>83.6</b>	1.45
<b>AM (MSS)</b>	<b>79.4</b>	<b>94.7</b>	<b>59.7</b>	83.5	<b>1.81</b>

TABLE IV

CONFIGURATIONS OF BOOSTING LOSS BY OSS ON FLICKR30K. \* DENOTES THE DBL STRATEGY WITH SOFT  $\gamma^{SA}$ . BOTH MODEL #H AND MODEL #I AVERAGE TWO BOOSTING LOSS DURING TRAINING PROCESS.

Model	Strategy				Sentence Retrieval		Image Retrieval		MD
	RS	RM	AS	AM	R@1	R@5	R@1	R@5	
A					75.8	93.4	56.5	81.4	0.91
B	✓				77.6	94.6	58.1	82.4	1.04
C			✓		78.2	94.8	57.1	82.5	1.13
D		✓*			77.1	94.5	<b>59.6</b>	<b>83.9</b>	1.23
E				✓*	77.2	94.6	58.8	83.5	1.30
F		✓			<b>79.3</b>	<b>95.4</b>	59.1	83.0	1.21
G				✓	79.0	94.9	58.5	83.1	<b>1.56</b>
H	✓		✓		78.0	94.6	57.8	82.3	1.13
I		✓		✓	78.9	95.1	58.7	83.0	1.44

within a mini-batch as a class and translates the cross-modal scalar projections into the normalized probability estimates between image and text. Besides, we employ  $H(P_s(x), P_t(x))$  to replace  $H(P_t(x), P_s(x))$  in Eq. (3) of DINO as DINO\* according to the observations [62], [80] that the latter could blur the distributions of multiple modes and bring ambiguities for image-text matching pattern. **1) Performance.** TABLE II shows that our DBL can consistently outweigh them on R@1 by a large margin (Maximum 2.9% and 1.7% on sentence and image retrieval), exhibiting good flexibility and broad applicability. Beyond the above, it also indicates that the knowledge transfers in a boosting manner hold tremendous potential for passing messages across branches and significantly promote the matching ability of single-branch learning. **2) Efficiency.** OAS/MSS requires no gradients for the anchor branch, and obtains about 12/11% memory and 120/18% time increase. Only OSS trains two branches simultaneously with extra 100% memory and 73% time costs. Notably, MSS can achieve both slight training time and memory costs. **3) Expansion.** Two-branch structure in OAS and MSS is a standard architecture of conventional distillation and contrastive learning. In OSS, DML [25] demonstrates another paradigm of multi-branch collaboration for mutual learning. Following it, we extend our strategies to larger peer cohorts. Concretely, we randomly initialize all branches  $\{B_i | i = 1, \dots, M\}$  where each branch  $B_m$  is supervised by regarding all previous cohorts  $\{B_i | i = 1, \dots, m - 1\}$  as the anchor branches. To ensure comparability between task-specific and boosting losses for current target branch, we average all boosting losses from its corresponding anchor branches, which are then added to the task-specific loss as the final objective function. We discover that the evaluation results slowly converge with the increasing number of peer branches and 3-branch DBL with RM earns the best

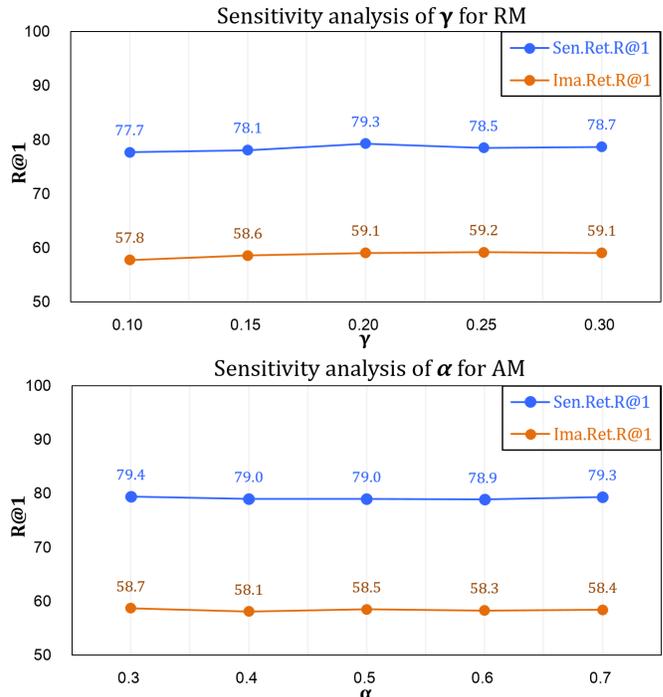


Fig. 4. Analyses of hyperparameter  $\gamma$  for RM and  $\alpha$  for AM under OSS.

bidirectional R@1 by 79.4/59.5% vs. 78.3/58.8% of 4-branch DML (best). The combination of multiple branches is flexible and two branches can balance complexity and performance.

**Comparison with Deep Metric Learning.** In TABLE III, we report the retrieval results with several popular loss functions based on single branch, which focus on cross-modal projection [62], hard negative mining [17], [18], and adaptive margin setting [20], [21] respectively. Here, we utilize our proposed DBL under MSS for comparison, given its relatively minimal additional consumption costs (only extra 11% memory and 18% time costs), when compared to the single branch with various loss functions. We can find that almost all similarity metrics learned through the above losses are better than the original ranking loss [15], confirming that it is beneficial to highlight the informative triplets and develop the appropriate thresholds. It is obvious that our DBL strategy achieves more impressive improvements and outperforms the best competitor MPL [18] on bidirectional R@1 by 1.3/1.7% and 2.3/2.3% in the relative/absolute manner separately. Besides, the comparison with the most related metrics [20], [21] validates that explicit margin constraints by the DBL are more effective and applicable for the network to obtain more powerful matching capabilities across modalities.

### C. Ablation Studies

**Configurations of boosting loss.** TABLE IV shows the different configurations of our DBL, consisting of Sum-Max sampling, Fixed-Soft margin value, and Relative-Absolute boosting manners. **1) Sum vs. Max.** Compared with Model #B and #C via Sum operation, Model #F and #G show that Max operation can excavate more valuable triplets and obtain the 1.7/0.8% and 1.0/1.4% R@1 increase on sentence and

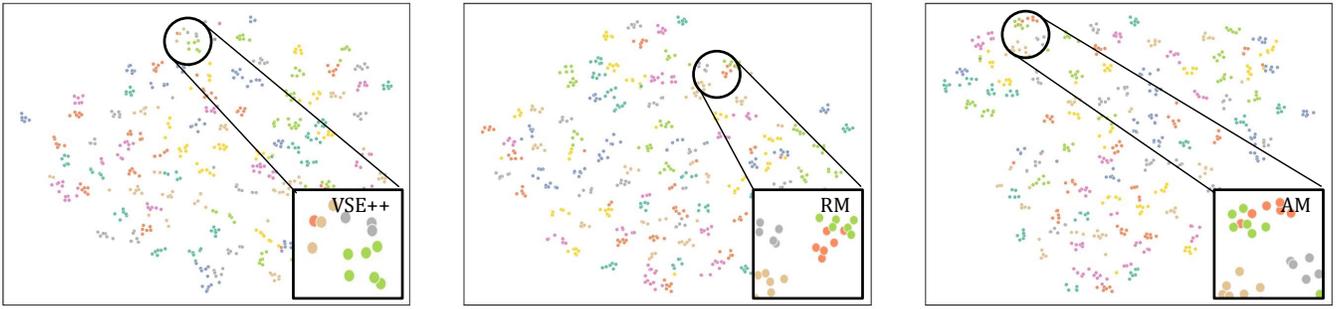


Fig. 5. Comparison of feature distribution for image and text samples. We implement t-SNE to visualize the image and text features based on GPO [42].

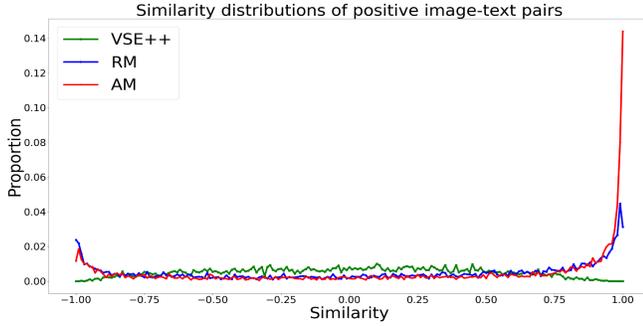


Fig. 6. Comparison of similarity distribution for positive image-text pairs.

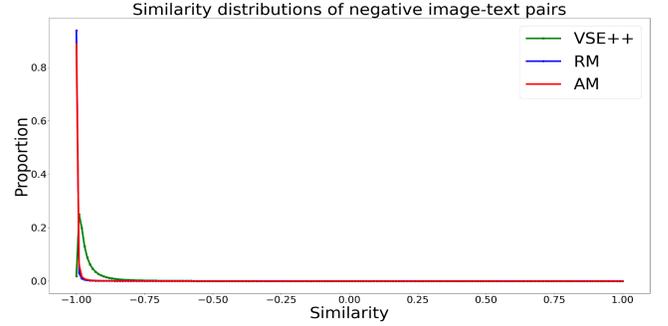


Fig. 7. Comparison of similarity distribution for negative image-text pairs.

image retrieval for the Relative/Absolute strategy. **2) Fixed vs. Soft.** Model #D, #E, #F, and #G demonstrate that the Soft form displays stronger abilities of image retrieval but achieves slightly attractive promotions on sentence retrieval, while the Fixed form obtains the optimal balance between image and sentence retrievals. Note that the Soft strategy stresses more on the hard triplets, making its gradient easily dominated by noise, being a result of either deficiency of model or data per se. **3) Relative vs. Absolute.** Comparing Model #B-#G, and more applications in TABLE I, we discover that both of them display the general effectiveness and unique qualities on  $R@k$  under diverse network architectures. Based on a more intuitive MD metric, the Absolute manner produces more compact similarity distributions than the Relative one. **4) Relative & Absolute.** We average the relative and absolute boosting losses to jointly supervise the training process. Compared with Model #F and #G, Model #H and #I display no further improvements on both  $R@k$  and MD metrics. This may be because AM imposes tighter constraints than RM, and such simple combinations fail to reinforce the penalty and further widen the distance between positive and negative pairs.

**Analyses of hyperparameter tuning.** The exclusive  $\alpha$  and  $\gamma$  tuning are shown in Fig. 4. For  $\alpha = 0.3 - 0.7$ , AM brings at least 3.2/1.9%  $R@1$  gains on sentence and image retrievals. For  $\gamma = 0.1 - 0.3$ , RM obtains a steady  $R@1$  gain of at least 1.9/1.3% at two directions. Note that the VSE++ loss [17] produces varied results with different margin  $\gamma$ , and  $\gamma = 0.2$  is a commonly-used configuration [36], [42], [51], [53]. Hence for simplicity and fair comparison, we set  $\alpha = 0.5$  and  $\gamma = 0.2$  in all our experiments, and obtain robust and stable performance benefits over various datasets and approaches.

**Impact of the well-trained anchor branch.** To validate this, we take AM under OAS on Flickr30K as an example. **(1)** We utilize  $\gamma=0.1$  of hard ranking loss [17] and 50% Flickr30k training data to construct two kinds of sub-optimal anchor branches. Interestingly, DBL still improves corresponding target branches by 2.7/1.9% and 7.4/5.8%  $R@1$  gains. **(2)** We also use VSE++ [17] as a poor backbone that outputs many noisy matching results, which achieves  $R@1$  benefits of 1.8/1.3% by DBL. **(3)** They verify that even with an inferior reference, DBL still displays good stability and robustness.

**Feature distributions of image and text samples.** We visualize the feature distribution of image and text features in Fig. 5. Note that for cross-modal interaction methods, feature visualization does not provide a direct reflection of the similarity measurement across modalities. Therefore, we adopt GPO [42] based on mono-modal representation for better illustration. We can observe that our DBL achieves better feature separation, even for some challenging samples. The anchor branch provides explicit and targeted distance information, which enables the target branch to enlarge the gap between image and text through our DBL strategy.

**Similarity distributions of image-text pairs.** Fig. 6 and 7 exhibit the similarity distributions of positive and negative pairs respectively on the Flickr30K test set. With VSE++ loss [17], the scores of negative pairs are concentrated near the value -1, while the curve of positive ones is relatively smooth. After introducing RM, their values are more densely distributed at the value -1 and 1 respectively, and the separability between them is fully exploited. Besides, AM can further enlarge the variations between matched and unmatched pairs, confirming that learning an adaptive and explicit margin



Fig. 8. Several retrieval examples on image retrieval. Green denotes the ground-truth image candidates and red denotes the unmatched retrieval samples.



Fig. 9. Several retrieval examples on sentence retrieval. Green denotes the ground-truth sentence candidates and red denotes the unmatched retrieval samples.

can lead to sufficient distance metrics and powerful matching patterns. Note that a slight peak arises around -1.0 for positive pairs. This is because DBL generally produces more confident predictions as compared to plain VSE++, and could inadvertently misclassify some hard positive pairs with weak correlations as negative samples. We further discover that the positive pairs with the similarity between -1.0 and -0.8 by RM and AM are mostly distributed ranging from -0.9 to -0.2 in the original VSE++ with 87.8% and 92.6% probability.

**Qualitative results of image and sentence retrieval.** Fig. 8 and 9 display several retrieval examples on image and

sentence retrieval, which can qualitatively indicate the learned distance measure between image and text. Compared with the anchor branch, the target branch by our DBL is capable of better recognizing cross-modal contents, and effectively distinguishing the correct results from various distractions with similar semantics, which validates the powerful applicability and matching capability of our DBL strategy.

V. CONCLUSION AND FUTURE WORKS

In this paper, we propose a novel Deep Boosting Learning (DBL) strategy to seek a powerful modeling capability by

imposing an adaptive and dynamic boosting mechanism for image-text matching task. Specifically, we first plumb the model property and data representation thoroughly, which in turn facilitates the learning process with appropriate regulations in a boosting manner. Extensive experiments on two benchmark datasets validate that our DBL further enlarges the insufficient variations within triplets and exploits the optimal feature and distance metrics across modalities. Besides, we discover that DBL consistently improves various popular frameworks by a large margin, confirming its general effectiveness and flexible applicability in image-text matching field. As one of the peer-to-peer strategies, DBL goes deeper than many related cooperative methods by learning margin knowledge to gain greater benefits under the same training and inference schemes, from which the community may get inspiration. In the future, we would like to incorporate them simultaneously in self-branch boosting and cross-branch imitating manners.

## REFERENCES

- [1] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, M. Ranzato, and T. Mikolov, "Devise: A deep visual-semantic embedding model," in *NIPS*, 2013, pp. 2121–2129. 1, 6
- [2] H. Wu, J. Mao, Y. Zhang, Y. Jiang, L. Li, W. Sun, and W. Ma, "Unified visual-semantic embeddings: Bridging vision and language with structured meaning representations," in *CVPR*, 2019, pp. 6609–6618. 1
- [3] S. Chen, Y. Zhao, Q. Jin, and Q. Wu, "Fine-grained video-text retrieval with hierarchical graph reasoning," in *CVPR*, 2020, pp. 10635–10644. 1
- [4] 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. 1
- [5] 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. 11895–11905. 1
- [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. 1
- [7] C. Gao, J. Chen, S. Liu, L. Wang, Q. Zhang, and Q. Wu, "Room-and-object aware knowledge reasoning for remote embodied referring expression," in *CVPR*, 2021, pp. 3064–3073. 1
- [8] Z. Yu, J. Yu, Y. Cui, D. Tao, and Q. Tian, "Deep modular co-attention networks for visual question answering," in *CVPR*, 2019, pp. 6281–6290. 1
- [9] X. Lin and D. Parikh, "Leveraging visual question answering for image-caption ranking," in *ECCV*, vol. 9906, 2016, pp. 261–277. 1
- [10] 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. 1
- [11] Q. Zhang, Z. Lei, Z. Zhang, and S. Z. Li, "Context-aware attention network for image-text retrieval," in *CVPR*, 2020, pp. 3533–3542. 1
- [12] 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. 1
- [13] 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. 1
- [14] J. Li, L. Niu, and L. Zhang, "Action-aware embedding enhancement for image-text retrieval," in *AAAI*, 2022, pp. 1323–1331. 1
- [15] L. Wang, Y. Li, and S. Lazebnik, "Learning deep structure-preserving image-text embeddings," in *CVPR*, 2016, pp. 5005–5013. 1, 2, 8
- [16] H. Nam, J. Ha, and J. Kim, "Dual attention networks for multimodal reasoning and matching," in *CVPR*, 2017, pp. 2156–2164. 1, 3
- [17] F. Faghri, D. J. Fleet, J. R. Kiros, and S. Fidler, "VSE++: improving visual-semantic embeddings with hard negatives," in *BMVC*, 2018, p. 12. 1, 2, 3, 6, 7, 8, 9
- [18] J. Wei, X. Xu, Y. Yang, Y. Ji, Z. Wang, and H. T. Shen, "Universal weighting metric learning for cross-modal matching," in *CVPR*, 2020, pp. 13002–13011. 1, 4, 8
- [19] X. Zhao, H. Qi, R. Luo, and L. Davis, "A weakly supervised adaptive triplet loss for deep metric learning," in *ICCV*, 2019, pp. 3177–3180. 1
- [20] A. F. Biten, A. Mafla, L. Gómez, and D. Karatzas, "Is an image worth five sentences? A new look into semantics for image-text matching," in *WACV*, 2022, pp. 2483–2492. 1, 2, 8
- [21] M. Zhou, Z. Niu, L. Wang, Z. Gao, Q. Zhang, and G. Hua, "Ladder loss for coherent visual-semantic embedding," in *AAAI*, 2020, pp. 13050–13057. 1, 2, 8
- [22] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "Fitnets: Hints for thin deep nets," in *ICLR*, 2015. 2
- [23] S. Zagoruyko and N. Komodakis, "Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer," in *ICLR*, 2017. 2
- [24] J. Yim, D. Joo, J. Bae, and J. Kim, "A gift from knowledge distillation: Fast optimization, network minimization and transfer learning," in *CVPR*, 2017, pp. 7130–7138. 2
- [25] Y. Zhang, T. Xiang, T. M. Hospedales, and H. Lu, "Deep mutual learning," in *CVPR*, 2018, pp. 4320–4328. 2, 5, 6, 7, 8
- [26] K. He, H. Fan, Y. Wu, S. Xie, and R. B. Girshick, "Momentum contrast for unsupervised visual representation learning," in *CVPR*, 2020, pp. 9726–9735. 2, 5, 6
- [27] J. Grill, F. Strub, F. Althé, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. Á. Pires, Z. Guo, M. G. Azar, B. Piot, K. Kavukcuoglu, R. Munos, and M. Valko, "Bootstrap your own latent - A new approach to self-supervised learning," in *NeurIPS*, 2020. 2, 5, 6
- [28] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, "Emerging properties in self-supervised vision transformers," in *ICCV*, 2021, pp. 9630–9640. 2, 5, 6, 7
- [29] G. E. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv: 1503.02531*, 2015. 2, 5
- [30] Y. Chen, N. Wang, and Z. Zhang, "Darkrank: Accelerating deep metric learning via cross sample similarities transfer," in *AAAI*, 2018, pp. 2852–2859. 2, 5, 6, 7
- [31] F. Tung and G. Mori, "Similarity-preserving knowledge distillation," in *ICCV*, 2019, pp. 1365–1374. 2, 5
- [32] W. Park, D. Kim, Y. Lu, and M. Cho, "Relational knowledge distillation," in *CVPR*, 2019, pp. 3967–3976. 2, 5, 6, 7
- [33] L. Yu, V. O. Yazici, X. Liu, J. van de Weijer, Y. Cheng, and A. Ramisa, "Learning metrics from teachers: Compact networks for image embedding," in *CVPR*, 2019, pp. 2907–2916. 2
- [34] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton, "A simple framework for contrastive learning of visual representations," in *ICML*, vol. 119, 2020, pp. 1597–1607. 2
- [35] X. Chen and K. He, "Exploring simple siamese representation learning," in *CVPR*, 2021, pp. 15750–15758. 2, 6
- [36] K. Li, Y. Zhang, K. Li, Y. Li, and Y. Fu, "Visual semantic reasoning for image-text matching," in *ICCV*, 2019, pp. 4653–4661. 2, 6, 7, 9
- [37] Y. Li, D. Zhang, and Y. Mu, "Visual-semantic matching by exploring high-order attention and distraction," in *CVPR*, 2020, pp. 12783–12792. 2
- [38] 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. 2
- [39] 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. 2
- [40] 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. 2
- [41] J. Li, L. Liu, L. Niu, and L. Zhang, "Memorize, associate and match: Embedding enhancement via fine-grained alignment for image-text retrieval," *IEEE TIP*, vol. 30, pp. 9193–9207, 2021. 2
- [42] J. Chen, H. Hu, H. Wu, Y. Jiang, and C. Wang, "Learning the best pooling strategy for visual semantic embedding," in *CVPR*, 2021, pp. 15789–15798. 2, 7, 9
- [43] H. Zhu, C. Zhang, Y. Wei, S. Huang, and Y. Zhao, "ESA: external space attention aggregation for image-text retrieval," *IEEE TCSVT*, vol. 33, no. 10, pp. 6131–6143, 2023. 2, 6, 7
- [44] 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. 2, 3, 6
- [45] 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. 2, 6, 7
- [46] 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. 2

- [47] 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. 2
- [48] 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. 2
- [49] H. Diao, Y. Zhang, W. Liu, X. Ruan, and H. Lu, "Plug-and-play regulators for image-text matching," *arXiv: 2303.13371*, 2023. 2
- [50] 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. 2, 6
- [51] H. Diao, Y. Zhang, L. Ma, and H. Lu, "Similarity reasoning and filtration for image-text matching," in *AAAI*, 2021, pp. 1218–1226. 2, 3, 6, 7, 9
- [52] 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. 2, 6
- [53] 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. 2, 6, 7, 9
- [54] A. Hermans, L. Beyer, and B. Leibe, "In defense of the triplet loss for person re-identification," *arXiv: 1703.07737*, 2017. 2, 4
- [55] W. Chen, X. Chen, J. Zhang, and K. Huang, "Beyond triplet loss: A deep quadruplet network for person re-identification," in *CVPR*, 2017, pp. 1320–1329. 2, 4
- [56] H. O. Song, Y. Xiang, S. Jegelka, and S. Savarese, "Deep metric learning via lifted structured feature embedding," in *CVPR*, 2016, pp. 4004–4012. 2
- [57] K. Sohn, "Improved deep metric learning with multi-class n-pair loss objective," in *NeurIPS*, 2016, pp. 1849–1857. 2
- [58] E. Ustinova and V. S. Lempitsky, "Learning deep embeddings with histogram loss," in *NeurIPS*, 2016, pp. 4170–4178. 2
- [59] Y. Movshovitz-Attias, A. Toshev, T. K. Leung, S. Ioffe, and S. Singh, "No fuss distance metric learning using proxies," in *ICCV*, 2017, pp. 360–368. 2
- [60] L. A. Hendricks, O. Wang, E. Shechtman, J. Sivic, T. Darrell, and B. C. Russell, "Localizing moments in video with natural language," in *ICCV*, 2017, pp. 5804–5813. 2
- [61] S. Li, T. Xiao, H. Li, B. Zhou, D. Yue, and X. Wang, "Person search with natural language description," in *CVPR*, 2017, pp. 5187–5196. 2
- [62] Y. Zhang and H. Lu, "Deep cross-modal projection learning for image-text matching," in *ECCV*, vol. 11205, 2018, pp. 707–723. 2, 7, 8
- [63] S. H. Lee, D. H. Kim, and B. C. Song, "Self-supervised knowledge distillation using singular value decomposition," in *ECCV*, vol. 11210, 2018, pp. 339–354. 2
- [64] 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. 3, 6
- [65] 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. 3
- [66] M. Schuster and K. K. Paliwal, "Bidirectional recurrent neural networks," *TSP*, vol. 45, no. 11, pp. 2673–2681, 1997. 3
- [67] P. F. Felzenszwalb, R. B. Girshick, D. A. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *TPAMI*, vol. 32, no. 9, pp. 1627–1645, 2010. 4
- [68] T. Malisiewicz, A. Gupta, and A. A. Efros, "Ensemble of exemplar-svms for object detection and beyond," in *ICCV*, 2011, pp. 89–96. 4
- [69] T. Chen, J. Deng, and J. Luo, "Adaptive offline quintuplet loss for image-text matching," in *ECCV*, vol. 12358, 2020, pp. 549–565. 4
- [70] N. Passalis, M. Tzelepi, and A. Tefas, "Probabilistic knowledge transfer for lightweight deep representation learning," *TNNLS*, vol. 32, no. 5, pp. 2030–2039, 2021. 5
- [71] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, "Unsupervised feature learning via non-parametric instance discrimination," in *CVPR*, 2018, pp. 3733–3742. 6
- [72] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," in *ECCV*, vol. 11218, 2018, pp. 139–156. 6
- [73] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, "Unsupervised learning of visual features by contrasting cluster assignments," in *NeurIPS*, 2020. 6
- [74] B. T. Polyak and A. B. Juditsky, "Acceleration of stochastic approximation by averaging," *SIAM journal on control and optimization*, vol. 30, no. 4, pp. 838–855, 1992. 6
- [75] D. Ruppert, "Efficient estimations from a slowly convergent robbins-monro process," Cornell University Operations Research and Industrial Engineering, Tech. Rep., 1988. 6
- [76] 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. 6
- [77] 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. 6
- [78] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *ICML*, vol. 139, 2021, pp. 8748–8763. 6, 7
- [79] X. Li, X. Yin, C. Li, P. Zhang, X. Hu, L. Zhang, L. Wang, H. Hu, L. Dong, F. Wei, Y. Choi, and J. Gao, "Oscar: Object-semantics aligned pre-training for vision-language tasks," in *ECCV*, vol. 12375, 2020, pp. 121–137. 6, 7
- [80] I. J. Goodfellow, Y. Bengio, and A. C. Courville, *Deep Learning*, ser. Adaptive computation and machine learning. MIT Press, 2016. 8