

Noisy-Correspondence Learning for Text-to-Image Person Re-identification

Yang Qin¹ Yingke Chen² Dezhong Peng^{1,4,5} Xi Peng¹ Joey Tianyi Zhou³ Peng Hu^{1*}

¹College of Computer Science, Sichuan University, Chengdu, 610095, China.

²Department of Computer and Information Sciences, Northumbria University, Newcastle upon Tyne NE1 8ST, UK.

³Centre for Frontier AI Research (CFAR) and Institute of High Performance Computing (IHPC), A*STAR, Singapore.

⁴Sichuan Newstrong UHD Video Technology Co., Ltd., Chengdu 610095, China.

⁵Chengdu Ruibei Yingte Information Technology Company Ltd., Chengdu 610065, China.

Abstract

Text-to-image person re-identification (TIReID) is a compelling topic in the cross-modal community, which aims to retrieve the target person based on a textual query. Although numerous TIReID methods have been proposed and achieved promising performance, they implicitly assume the training image-text pairs are correctly aligned, which is not always the case in real-world scenarios. In practice, the image-text pairs inevitably exist under-correlated or even false-correlated, a.k.a noisy correspondence (NC), due to the low quality of the images and annotation errors. To address this problem, we propose a novel Robust Dual Embedding method (RDE) that can learn robust visual-semantic associations even with NC. Specifically, RDE consists of two main components: 1) A Confident Consensus Division (CCD) module that leverages the dual-grained decisions of dual embedding modules to obtain a consensus set of clean training data, which enables the model to learn correct and reliable visual-semantic associations. 2) A Triplet Alignment Loss (TAL) relaxes the conventional Triplet Ranking loss with the hardest negative samples to a log-exponential upper bound over all negative ones, thus preventing the model collapse under NC and can also focus on hard-negative samples for promising performance. We conduct extensive experiments on three public benchmarks, namely CUHK-PEDES, ICFG-PEDES, and RSTPReID, to evaluate the performance and robustness of our RDE. Our method achieves state-of-the-art results both with and without synthetic noisy correspondences on all three datasets. Code is available at <https://github.com/QinYang79/RDE>.

1. Introduction

Text-to-image person re-identification (TIReID) [24, 27, 45] aims to understand the natural language descriptions

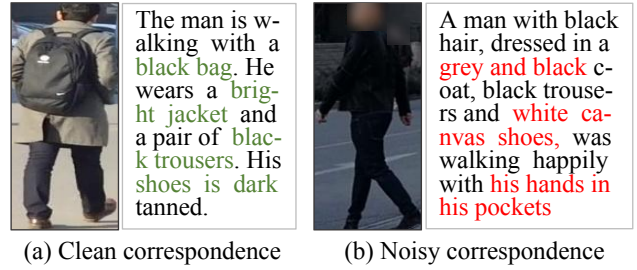


Figure 1. The illustration of noisy correspondence. The figure shows an example of the NC problem, which occurs when the image-text pairs are wrongly aligned, i.e., false positive pairs (FPPs). Since the model does not know which pairs are noisy in practice, they will unavoidably degrade the performance by incorrect supervision information. As seen in the figure, (a) the clean image-text pair is semantically matched, while (b) the noisy pair is not, which would cause the cross-modal model to learn erroneous visual-textual associations. Note that both examples in (a) and (b) are from and actually exist in the RSTPReid dataset [62].

to retrieve the matched person image from a large gallery set. This task has received increasing attention from both academic and industrial communities recently, e.g., finding/tracking suspect/lost persons in a surveillance system [11, 47]. However, TIReID remains a challenging task due to the inherent heterogeneity gap across different modalities and appearance attribute redundancy.

To tackle these challenges, most of the existing methods explore global- and local-matching alignment to learn accurate similarity measurements for person re-identification. To be specific, some global-matching methods [45, 52, 60] leverage vision/language backbones to extract modality-specific features and employ contrastive learning to achieve global visual-semantic alignments. To capture fine-grained information, some local-matching methods [25, 34, 42, 46] explicitly align local body regions to textually described entities/objectives to improve the discriminability of pedestrian features. Recently, some works [19, 24, 53] pro-

*Corresponding author: Peng Hu (penghu.ml@gmail.com).

pose to exploit visual/semantic knowledge learned by the pre-trained models, such as BERT [6], ViT [10], and CLIP [40], and achieve explicit global alignments or discover more fine-grained local correspondence, thus boosting the re-identification performance. Although these methods achieve remarkable progress, they implicitly assume that all training image-text pairs are aligned correctly.

In reality, this assumption is hard or even impossible to hold due to the person’s pose, camera angle, illumination, and other inevitable factors in images, which may result in some inaccurate/mismatched textual descriptions of images (see Figure 1), *e.g.*, the RSTPReid dataset [62]. Moreover, we observe that excessive such imperfect/mismatched image-text pairs would cause an overfitting problem and degrade the performance of existing TIREID methods shown in Figure 5. Based on the observation, in this paper, we reveal and study a new problem in TIREID, *i.e.*, noisy correspondence (NC). Different from noisy labels, NC refers to the false correspondences of image-text pairs in TIREID, *i.e.*, False Positive Pairs (FPPs): some negative image-text pairs are used as positive ones for cross-modal learning. Inevitably, FPPs will misguide models to overfit noisy supervision and collapse to suboptimal solutions due to the memorization effect [1] of Deep Neural Networks (DNNs).

To address the NC problem, we propose a Robust Dual Embedding method (RDE) for TIREID in this paper, which benefits from an effective Confident Consensus Division mechanism (CCD) and a novel Triplet Alignment Loss (TAL). Specifically, CCD fuses the dual-grained decisions to consensually divide the training data into clean and noisy sets, thus providing more reliable correspondences for robust learning. To diversify the model grain, the basic global embeddings (BGE) and token selection embeddings (TSE) are presented for coarse-grained and fine-grained cross-modal interactions respectively, thus capturing visual-semantic associations comprehensively. Different from the widely-used Triplet Ranking loss with the hardest negatives, our TAL relaxes the similarity learning from the hardest negative samples to all negative ones by applying an upper bound, which brings a stable solution for the collapse of training under NC while also benefiting from the hardest negatives mining to achieve promising performance. As a result, our RDE can achieve robustness against NC thanks to the proposed reliable supervision and stable triplet loss. The contributions and innovations of this paper are summarized as follows:

- We reveal and study a new and ubiquitous problem in TIREID, termed noisy correspondence (NC). Different from class-level noisy labels, NC refers to erroneous correspondences in the person-description pairs that can mislead the model to learn incorrect visual-semantic associations. To the best of our knowledge, this paper could be the first work to explore this problem in TIREID.

- We propose a robust method, termed RDE, to mitigate the adverse impact of NC through the proposed Confident Consensus Division (CCD) and novel Triplet Alignment Loss (TAL). By using CCD and TAL, RDE can obtain convincing consensus pairs and reduce the misleading risks in training, thus embracing robustness against NC.
- Extensive experiments on three public image-text person benchmarks demonstrate the robustness and superiority of our method. Our method achieves the best performance both with and without synthetic noisy correspondence on all three datasets.

2. Related Work

2.1. Text-to-Image Person Re-identification

Text-to-image person re-identification (TIREID) is a novel and challenging task that aims to match a person image with a given natural language description [2–4, 27, 29, 33, 43, 44, 55, 60]. Existing TIREID methods could be roughly classified into two groups according to their alignment levels, *i.e.*, **global-matching** methods [45, 61, 62] and **local-matching** methods [16, 42, 46]. The former try to learn cross-modal embeddings in a common latent space by employing textual and visual backbones with a matching loss (*e.g.*, CPM/C loss [60] and Triplet Ranking loss [12]) for TIREID. However, these methods mainly focus on global features while ignoring the fine-grained interactions between local features, which limits their performance improvement. To achieve fine-grained interactions, some of the latter methods explore explicit local alignments between body regions and textual entities for more refined alignments. However, these methods require more computational resources due to the complex local-level associations. Recently, inspired and benefited from **vision-language pre-training** models [40], some methods [19, 24, 53] expect to use the learned rich alignment knowledge of pre-trained models for local- or global-alignments. Although these methods achieve promising performance, almost all of them implicitly assume that all input training pairs are correctly aligned, which is hard to meet in practice due to the ubiquitous noise. In this paper, we address the inevitable and challenging noisy correspondence problem in TIREID.

2.2. Learning with Noisy Correspondence

As a special learning paradigm with noisy labels [14, 26, 32] in multi-modal/view community [21, 37, 37, 38, 57], the studies for noisy correspondence (NC) have recently attracted more and more attention in various tasks, *e.g.*, video-text retrieval [59], visible-infrared person re-identification [31, 56, 58], and image-text matching [23, 36], which means that the negative pairs are wrongly treated as positive ones, *i.e.*, false positive pairs (FPPs). To handle

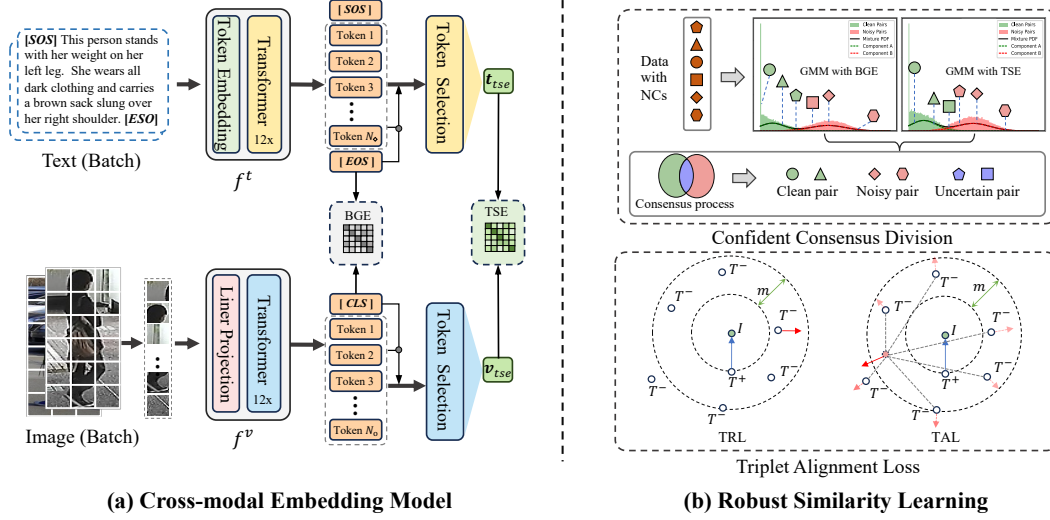


Figure 2. The overview of our RDE. (a) is the illustration of the cross-modal embedding model used in RDE, which consists of *basical global embedding* (BGE) and *token selection embedding* (TSE) modules with different granularity. By integrating them, RDE can capture coarse-grained cross-modal interactions while selecting informative local token features to encode more fine-grained representations for a more accurate similarity. (b) shows the core of RDE to achieve robust similarity learning, which consists of Confident Consensus Division (CCD) and Triplet Alignment Loss (TAL). CCD performs consensus division to obtain confident clean training data, thus avoiding misleading from noisy pairs. Unlike traditional Triplet Ranking Loss (TRL) [12], TAL exploits an upper bound to consider all negative pairs, thus embracing more stable learning.

this problem, numerous methods are proposed to learn with NC, which can be broadly categorized into **sample selection** [18, 23, 59] and **robust loss functions** [22, 36, 39, 56]. The former commonly leverage the memorization effect of DNNs [1] to gradually distinguish the noisy data, thus paying more attention to clean data while less attention to noisy data. Differently, the latter methods aim to develop noise-tolerance loss functions to improve the robustness of model training against NC. Although the aforementioned methods achieve promising performance in various tasks, they are not specifically designed for TIReID and may be inefficient or ineffective in person re-identification. In this paper, we propose a well-designed method to tackle the NC problem in TIReID, which not only performs superior in noisy scenarios but also achieves promising performance in ordinary scenarios.

3. Methodology

3.1. Problem Statement

The purpose of TIReID is to retrieve a pedestrian image from the gallery set that matches the given textual description. For clarity, we represent the gallery set as $\mathcal{V} = \{I_i, y_i^p, y_i^v\}_{i=1}^{N_v}$ and the corresponding text set as $\mathcal{T} = \{T_i, y_i^p\}_{i=1}^{N_t}$, where N_v is the number of images, N_t is the number of texts, $y_i^p \in \mathcal{Y}_p = \{1, \dots, C\}$ is the class label (person identify), C is the number of identifies, and $y_i^v \in \mathcal{Y}_v = \{1, \dots, N_v\}$ is the image label. The image-text pair set used in TIReID can be defined as

$\mathcal{P} = \{(I_i, T_i), y_i^v, y_i^p\}_{i=1}^N$, where the cross-modal samples of each pair have the same image label y_i^v and class label y_i^p . We define a binary correspondence label $l_{ij} \in \{0, 1\}$ to indicate the matched degree of any image-text pair. If $l_{ij} = 1$, the pair (I_i, T_j) is matched (positive pair), otherwise it is not (negative pair). In practice, due to ubiquitous annotation noise, some unmatched pairs ($l_{ij} = 0$) are wrongly labeled as matched ($l_{ij} = 1$), resulting in noisy correspondences (NCs) and performance degradation. To handle NC for robust TIReID, we present an RDE that leverages the Confident Consensus Division (CCD) and Triplet Alignment Loss (TAL) to mitigate the negative impact of label noise.

3.2. Cross-modal Embedding Model

In this section, we describe the cross-modal model used in our RDE. Following previous work [24], we utilize the visual encoder f^v and textual encoder f^t of the pre-trained model CLIP as modality-specific encoders to obtain token representations and implement cross-modal interactions through two embedding modules.

3.2.1 Token Representations

Give an input image $I_i \in \mathcal{V}$, we use the visual encoder f^v of CLIP to tokenize the image into a discrete token representation sequence with a length of $N_o + 1$, i.e., $\mathbf{V}_i = f^v(I_i) = \{\mathbf{v}_g^i, \mathbf{v}_1^i, \mathbf{v}_2^i, \dots, \mathbf{v}_{N_o}^i\}^\top \in \mathbb{R}^{(N_o+1) \times d}$, where d is the dimensionality of the shared latent space. These fea-

tures include an encoded feature \mathbf{v}_g^i of the [CLS] token and patch-level local features $\{\mathbf{v}_j^i\}_{j=1}^{N_\diamond}$ of N_\diamond fixed-sized non-overlapping patches of I_i , wherein \mathbf{v}_g^i can represent the global representation. For an input text $T_i \in \mathcal{T}$, we apply the textual encoder f^t of CLIP to obtain global and local representations. Specifically, following IRRRA [24], we first tokenize the input text T_i using lower-cased byte pair encoding (BPE) with a 49,152 vocab size into a token sequence. The token sequence is bracketed with [SOS] and [EOS] tokens to represent the beginning and end of the sequence. Then, we feed the token sequence into f_t to obtain the features $\mathbf{T}_i = \{\mathbf{t}_s^i, \mathbf{t}_1^i, \dots, \mathbf{t}_{N_\diamond}^i, \mathbf{t}_e^i\}^\top \in \mathbb{R}^{(N_\diamond+2) \times d}$, where \mathbf{t}_s^i and \mathbf{t}_e^i are the features of [SOS] and [EOS] tokens and $\{\mathbf{v}_j^i\}_{j=1}^{N_\diamond}$ are the word-level local features of N_\diamond word tokens of text T_i . Generally, the \mathbf{t}_e^i can be regarded as the sentence-level global feature of T_i .

3.2.2 Dual Embedding Modules

To measure the similarity between any image-text pair (I_i, T_j) , we can directly use the global features of [CLS] and [EOS] tokens to compute the *Basic Global Embedding* (BGE) similarity by the cosine similarity, i.e., $S_{ij}^b = \mathbf{v}_g^i \mathbf{t}_e^j^\top / \|\mathbf{v}_g^i\| \|\mathbf{t}_e^j\|$, where the global features represent the global embedding representations of two modalities. However, optimizing the BGE similarities alone may not capture the fine-grained interactions between two modalities, which will limit performance improvement. To address this issue, we exploit the local features of informative tokens to learn more discriminative embedding representations, thus mining the fine-grained correspondences. In CLIP, the global features of the tokens ([CLS] and [EOS]) are obtained by a weighted aggregation of all local token features. These weights reflect the correlation between the global token and each local token. Following previous methods [53, 63], we could select the informative tokens based on these correlation weights to aggregate local features for a more representative global embedding.

In practice, these correlation weights can be obtained directly in the self-attention map of the last Transformer blocks of f^v and f^t , which reflects the relevance among the input $1 + N_\diamond$ (or $2 + N_\diamond$) tokens. Given the output self-attention map $\mathbf{A}_i^v \in \mathbb{R}^{(1+N_\diamond) \times (1+N_\diamond)}$ of image I_i , the correlation weights between global token and local tokens are $\{a_{i,j}^v\}_{j=1}^{N_\diamond} = \mathbf{a}_i^v = \mathbf{A}_i^v[0, 1 : N_\diamond + 1] \in \mathbb{R}^{N_\diamond}$. Similarly, for text T_i , the correlation weights are $\{a_{i,j}^t\}_{j=1}^{N_\diamond} = \mathbf{a}_i^t = \mathbf{A}_i^t[0, 1 : N_\diamond + 1] \in \mathbb{R}^{N_\diamond}$, where $\mathbf{A}_i^t \in \mathbb{R}^{(2+N_\diamond) \times (2+N_\diamond)}$ is the output self-attention map for text T_i . Then, we select a proportion (*TopK*) of the corresponding token features with higher scores for embedding. Specifically, for I_i , the selected token sequences and correlation weights are reorganized as $\mathbf{V}_i^s = \{\mathbf{v}_j^i\}_{j \in \mathbf{K}_i^v}$ and $\hat{\mathbf{a}}_i^v = \{a_{i,j}^v\}_{j \in \mathbf{K}_i^v}$, where

$\mathbf{K}_i^v \in \mathbb{R}^{\lfloor \mathcal{R} N_\diamond \rfloor}$ is the set of indices for the selected local tokens of I_i and \mathcal{R} is the selection ratio. For text T_i , the selected token sequences and correlation weights are also reorganized as $\mathbf{T}_i^s = \{\mathbf{t}_j^i\}_{j \in \mathbf{K}_i^t}$ and $\hat{\mathbf{a}}_i^t = \{a_{i,j}^t\}_{j \in \mathbf{K}_i^t}$, where $\mathbf{K}_i^t \in \mathbb{R}^{\min(\lfloor \mathcal{R} N_\diamond' \rfloor, N_\diamond)}$ is the set of indices for the selected local tokens of T_i . N_\diamond' is the maximum input sequence length of f^t . For I_i and T_i , we perform an embedding transformation on these selected token features to obtain subtle representations, instead of using complex fine-grained correspondence discovery used in CFine [53]. The transformation is performed by an embedding module like the residual block [20], as follows:

$$\begin{aligned} \mathbf{v}_{tse}^i &= \text{MaxPool}(\text{MLP}(\hat{\mathbf{V}}_i^s) + \text{FC}(\hat{\mathbf{V}}_i^s)), \\ \mathbf{t}_{tse}^i &= \text{MaxPool}(\text{MLP}(\hat{\mathbf{T}}_i^s) + \text{FC}(\hat{\mathbf{T}}_i^s)), \end{aligned} \quad (1)$$

where $\text{MaxPool}(\cdot)$ is the max-pooling function, $\text{MLP}(\cdot)$ is a multi-layer perceptron (MLP) layer, $\text{FC}(\cdot)$ is a linear layer, $\hat{\mathbf{V}}_i^s = \text{L2Norm}(\mathbf{V}_i^s)$, and $\hat{\mathbf{T}}_i^s = \text{L2Norm}(\mathbf{T}_i^s)$. $\text{L2Norm}(\cdot)$ is the ℓ_2 -normalization function to normalize features. Finally, for any pair (I_i, T_j) , we compute the cosine similarity S_{ij}^t between \mathbf{v}_{tse}^i and \mathbf{t}_{tse}^j as the *Token Selection Embedding* (TSE) similarity to measure the cross-modal matching degree for auxiliary training and inference.

3.3. Robust Similarity Learning

In this section, we detail how we use the image-text similarities computed by the dual embedding modules for robust TIReID, which involves Confident Consensus Division (CCD) and Triplet Alignment Loss (TAL).

3.3.1 Confident Consensus Division

To alleviate the negative impact of NC, the key is to filter the possible noisy pairs in the training data, which directly avoids false supervision information. Some previous work in learning with noisy labels [17, 23, 26] are inspired by the memorization effect [1] of DNNs to perform filtrations, i.e., the clean (easy) data tend to have a smaller loss value than that of noisy (hard) data in early training. Based on this, we can exploit the two-component Gaussian Mixture Model (GMM) to fit the per-sample loss distributions computed by the predictions of BGE and TSE to further identify the noisy pairs in the training data. Specifically, given a cross-modal model \mathcal{M} , we first define the per-sample loss as:

$$\ell(\mathcal{M}, \mathcal{P}) = \{\ell_i\}_{i=1}^N = \{\mathcal{L}(I_i, T_i)\}_{i=1}^N, \quad (2)$$

where \mathcal{L} is the loss function for pair $(I_i, T_i) \in \mathcal{P}$ to bring them closer in the shared latent space. In our method, \mathcal{L} is the proposed \mathcal{L}_{tal} defined in Equation (11). Then, the per-sample loss is fed into the GMM to separate clean and noisy data, i.e., assigning the Gaussian component with a

lower mean value as a clean set and the other as a noisy one, respectively. Following [23, 26], we use the Expectation-Maximization algorithm to optimize the GMM and compute the posterior probability $p(k|\ell_i) = p(k)p(\ell_i|k)/p(\ell_i)$ for the i -th pair as the probability of being clean/noisy pair, where $k \in \{0, 1\}$ is used to indicate whether it is a clean or a noisy component. Then, we set a threshold $\delta = 0.5$ to $\{p(k = 0|\ell_i)\}_{i=1}^N$ to divide the data into clean and noisy sets, *i.e.*,

$$\begin{aligned} \mathcal{P}^c &= \{(I_i, T_i) | p(k = 0|\ell_i) > \delta, \forall (I_i, T_i) \in \mathcal{P}\}, \\ \mathcal{P}^n &= \{(I_i, T_i) | p(k = 0|\ell_i) \leq \delta, \forall (I_i, T_i) \in \mathcal{P}\}, \end{aligned} \quad (3)$$

where \mathcal{P}^c and \mathcal{P}^n are the divided clean and noisy sets, respectively. For BGE and TSE, the divisions conducted with Equation (13) are $\mathcal{P} = \mathcal{P}_{bge}^c \cup \mathcal{P}_{bge}^n$ and $\mathcal{P} = \mathcal{P}_{tse}^c \cup \mathcal{P}_{tse}^n$, separately.

To obtain the final reliable divisions, we propose to exploit the consistency between the two divisions to find the consensus part as the final confident clean set, *i.e.*, $\hat{\mathcal{P}}^c = \mathcal{P}_{bge}^c \cap \mathcal{P}_{tse}^c$. The rest of the data can be divided into noisy and uncertain subsets, *i.e.*, $\hat{\mathcal{P}}^n = \mathcal{P}_{bge}^n \cap \mathcal{P}_{tse}^n$ and $\hat{\mathcal{P}}^u = \mathcal{P} - (\hat{\mathcal{P}}^c \cup \hat{\mathcal{P}}^n)$. Finally, we exploit the divisions to further recalibrate the correspondence labels, *e.g.*, for i -th pair, the process can be expressed as:

$$\hat{l}_{ii} = \begin{cases} 1, & \text{if } (I_i, T_i) \in \hat{\mathcal{P}}^c, \\ 0, & \text{if } (I_i, T_i) \in \hat{\mathcal{P}}^n, \\ \text{Rand}(\{0, 1\}), & \text{if } (I_i, T_i) \in \hat{\mathcal{P}}^u, \end{cases} \quad (4)$$

where $\text{Rand}(X)$ is the function to randomly select an element from the collection X .

3.3.2 Triplet Alignment Loss

The Triplet Ranking Loss (TRL) is a common matching loss that is widely used in cross-modal learning, and achieves promising performance by employing the hardest negatives, *e.g.*, image-text matching [7], video-text retrieval [9], *etc.* However, we find that this strategy may lead to bad local minima or even model collapse for TIReID under NC in the early stages of training. In contrast, the summation version of TRL that considers all negative samples, namely TRL-S, can maintain better stability and avoid model collapse, but suffers from insufficient performance due to the lack of attention to hard negatives (see Section 3.3.3 for more discussion). Therefore, we propose a novel Triplet Alignment Loss (TAL) to guide TIReID, which differs from TRL in that it relaxes the optimization of the hardest negatives to all negatives with an upper bound (see Lemma 1). Thanks to the relaxation, TAL reduces the risk of the optimization being dominated by the hardest negatives, thereby making the training more stable and comprehensive by considering

all pairs. For an input pair (I_i, T_i) in a mini-batch \mathbf{x} , TAL is defined as

$$\begin{aligned} \mathcal{L}_{tal}(I_i, T_i) &= [m - S_{i2t}^+(I_i) + \tau \log(\sum_{j=1}^K q_{ij} \exp(S(I_i, T_j)/\tau))]_+ \\ &\quad + [m - S_{i2t}^+(T_i) + \tau \log(\sum_{j=1}^K q_{ji} \exp(S(I_j, T_i)/\tau))]_+, \end{aligned} \quad (5)$$

where m is a positive margin coefficient, τ is a temperature coefficient to control hardness, $S(I_i, T_j) \in \{S_{ij}^b, S_{ij}^t\}$, $[x]_+ \equiv \max(x, 0)$, $\exp(x) \equiv e^x$, $q_{ij} = 1 - l_{ij}$, and K is the size of \mathbf{x} . From Lemma 1, as $\tau \rightarrow 0$, TAL approaches TRL and focuses more on hard negatives. Since multiple positive pairs from the same identity may appear in the mini-batch, $S_{i2t}^+(I_i) = \sum_{j=1}^K \alpha_{ij} S(I_i, T_j)$ is the weighted average similarity of positive pairs for image I_i , where $\alpha_{ij} = \frac{l_{ij} \exp(S(I_i, T_j)/\tau)}{\sum_{k=1}^N l_{ik} \exp(S(I_i, T_k)/\tau)}$. Similarly, $S_{i2t}^+(T_i)$ is the weighted average similarity of positive pairs for text T_i .

Lemma 1 TAL is the upper bound of TRL, *i.e.*,

$$\begin{aligned} \mathcal{L}_{trl}(I_i, T_i) &= [m - S_{i2t}^+(I_i) + S(I_i, \hat{T}_i)]_+ \\ &\quad + [m - S_{i2t}^+(T_i) + S(\hat{I}_i, T_i)]_+ \leq \mathcal{L}_{tal}(I_i, T_i), \end{aligned} \quad (6)$$

where $\hat{T}_i \in \{T_j | l_{ij} = 0, \forall j \in \{1, \dots, K\}\}$ is the hardest negative text for I_i and $\hat{I}_i \in \{I_j | l_{ji} = 0, \forall j \in \{1, \dots, K\}\}$ is the hardest negative image for I_i , respectively.

3.3.3 Revisit Triplet Raking Loss

To explore the behaviors of the triplet losses in the noisy case, we record the similarity distributions versus iterations of TRL, TRL-S, and the proposed TAL under 50% noise. From Figure 3a, one can see that the similarities of all pairs are gradually gathered to 1 during training with TRL, *i.e.*, all samples *collapses* to a narrow neighborhood space on a hypersphere, resulting in a trivial solution and a bad performance (3.64%). To delve deeper into the underlying rea-

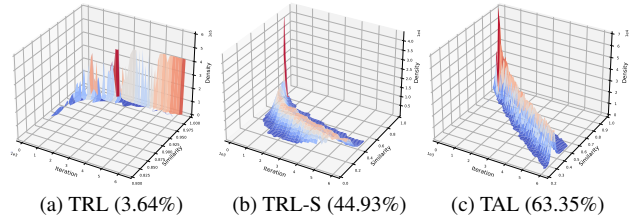


Figure 3. The difference between TRL, TRL-S, and proposed TAL on the similarity distribution versus iterations. The y - z plane represents the similarity density. The corresponding Rank-1 scores of testing are placed in brackets for convenience.

son, we performed a gradient analysis. For ease of representation and analysis, we only consider one direction since

image-to-text retrieval and text-to-image retrieval are symmetrical. And, we suppose that there is only one paired text for each image in the mini-batch. Due to the truncation operation $[x]_+$, we only discuss the case of $\mathcal{L} > 0$ that could generate gradients. Taking the image-to-text direction as an example, the gradients generated by TRL, TRL-S, and TAL are

$$\frac{\partial \mathcal{L}_{trl}}{\partial \mathbf{v}_i} = \hat{\mathbf{t}}_i - \mathbf{t}_i, \quad \frac{\partial \mathcal{L}_{trl}}{\partial \mathbf{t}_i} = -\mathbf{v}_i, \quad \frac{\partial \mathcal{L}_{trl}}{\partial \mathbf{t}_j} = \mathbf{v}_i, \quad (7)$$

$$\frac{\partial \mathcal{L}_{trls}}{\partial \mathbf{v}_i} = \sum_{j \in \mathcal{Z}} (\mathbf{t}_j - \mathbf{t}_i), \quad \frac{\partial \mathcal{L}_{trls}}{\partial \mathbf{t}_i} = -|\mathcal{Z}| \mathbf{v}_i, \quad \frac{\partial \mathcal{L}_{trls}}{\partial \mathbf{t}_j} = \mathbf{v}_i, \quad (8)$$

$$\frac{\partial \mathcal{L}_{tal}}{\partial \mathbf{v}_i} = \sum_{j \neq i}^K \beta_j (\mathbf{t}_j - \mathbf{t}_i), \quad \frac{\partial \mathcal{L}_{tal}}{\partial \mathbf{t}_i} = -\mathbf{v}_i, \quad \frac{\partial \mathcal{L}_{tal}}{\partial \mathbf{t}_j} = \beta_j \mathbf{v}_i, \quad (9)$$

where $\mathcal{Z} = \{z \mid [m - S(I_i, T_i) + S(I_i, T_z)]_+ > 0, z \neq i, z \in \{0, \dots, K\}\}$, $\beta_j = \frac{\exp(\mathbf{v}_i^\top \mathbf{t}_j / \tau)}{\sum_{k \neq i}^K \exp(\mathbf{v}_i^\top \mathbf{t}_k / \tau)}$, $\hat{\mathbf{t}}_i$, \mathbf{t}_j and \mathbf{t}_i are the hardest negative sample, negative sample, and positive sample of the anchor sample \mathbf{v}_i , respectively. Since the hardest sample is most similar to the positive one, $\frac{\partial \mathcal{L}_{trl}}{\partial \mathbf{v}_i}$ would easily approach 0 and the gradients for other negative samples except for the hardest negative one are all 0, which may lead to bad local minima early on in training and even cause the worst-case scenario, *i.e.*, model collapse (see Figure 3a). Unlike TRL, TRL-S aims to push all negative samples away from the anchor by a constant margin and produces stronger gradients for the anchor, *i.e.*, $\|\frac{\partial \mathcal{L}_{trls}}{\partial \mathbf{v}_i}\|_2 \geq \|\frac{\partial \mathcal{L}_{trl}}{\partial \mathbf{v}_i}\|_2$, thus avoiding model collapse (see Figure 3b). However, the drawback is that TRL-S treats every negative sample equally while ignoring challenging ones, which limits performance improvement. Different from TRL and TRL-S, from Equation (22), our TAL can comprehensively consider all negative samples and exploits the anchor-negative semantic relationships to adaptively adjust the gradients for each negative, thus paying more attention to hard negatives. As a result, TAL would avoid model collapse under NC while achieving superior performance (63.35% vs. 44.93% vs. 3.64%). More details for the derivations of gradients are provided in the supplementary material.

3.3.4 Training and Inference

To train the model robustly, we use the corrected label \hat{l}_{ii} instead of the original correspondence label l_{ii} to compute the final matching loss, *i.e.*,

$$\mathcal{L}_m = \sum_{i=1}^K \hat{l}_{ii} (\mathcal{L}^b(I_i, T_i) + \mathcal{L}^t(I_i, T_i)), \quad (10)$$

where $\mathcal{L}^b(I_i, T_i)$ and $\mathcal{L}^t(I_i, T_i)$ are the TAL losses computed by Equation (11) with BGE and TSE similarities, respectively. The training process of RDE is shown in Algorithm 1. For the joint inference, we compute the final

Algorithm 1 The training process of our RDE

Input: The training data \mathcal{P} with N image-text pairs, maximal epoch N_e , the cross-modal model $\mathcal{M}(\Theta)$, and the hyper-parameters \mathcal{R}, m, τ ;

- 1: Initialize the backbones with the weights of the pre-trained CLIP except for the TSE module, which is randomly initialized;
- 2: **for** $e = 1, 2, \dots, N_e$ **do**
- 3: Calculate the per-sample loss $\ell(\mathcal{M}, \mathcal{P})$;
- 4: Divide the training data with the predictions of BGE and TSE using Equation (13), respectively;
- 5: Obtain the consensus divisions to recalibrate the correspondence labels $\{\hat{l}_{ii}\}_{i=1}^N$ with Equation (4);
- 6: **for** \mathbf{x} in mini-batches $\{\mathbf{x}_m\}_{m=1}^M$ **do**
- 7: Extract the BGE and TSE features of \mathbf{x} ;
- 8: Compute the similarities between K image-text pairs in \mathbf{x} with above features;
- 9: Calculate the final matching loss \mathcal{L}_m with Equation (10);
- 10: $\Theta = \text{Optimizer}(\Theta, \mathcal{L}_m)$;
- 11: **end for**
- 12: **end for**

Output: The optimized parameters $\hat{\Theta}$.

similarity of the image-text pair as the average of the similarities computed by both embedding modules, *i.e.*, $S = (S^b + S^t)/2$.

4. Experiments

In this section, we conduct extensive experiments to verify the effectiveness and superiority of the proposed RDE on three widely-used benchmark datasets.

4.1. Datasets and Settings

4.1.1 Datasets

In the experiments, we use the CHUK-PEDES [27], ICFG-PEDES [8], and RSTPreid [62] datasets to evaluate our RDE. We follow the data partitions used in IRRA [24] to split the datasets into training, validation, and test sets, wherein the ICFG-PEDES dataset only has training and validation sets. More details are provided in the supplementary material.

4.1.2 Evaluation Protocols

For all experiments, we mainly employ the popular Rank-K metrics ($K=1, 5, 10$) to measure the retrieval performance. In addition to Rank-K, we also adopt the mean Average Precision (mAP) and mean Inverse Negative Penalty (mINP) as auxiliary retrieval metrics to further evaluate performance following [24].

Noise	Methods		CUHK-PEDES					ICFG-PEDES					RSTPReid				
			R-1	R-5	R-10	mAP	mINP	R-1	R-5	R-10	mAP	mINP	R-1	R-5	R-10	mAP	mINP
0%	SSAN	Best	61.37	80.15	86.73	-	-	54.23	72.63	79.53	-	-	43.50	67.80	77.15	-	-
	IVT	Best	65.59	83.11	89.21	-	-	56.04	73.60	80.22	-	-	46.70	70.00	78.80	-	-
	CFine	Best	69.57	85.93	91.15	-	-	60.83	76.55	82.42	-	-	50.55	72.50	81.60	-	-
	IRRA	Best	73.38	89.93	93.71	<u>66.13</u>	<u>50.24</u>	<u>63.46</u>	<u>80.25</u>	<u>85.82</u>	<u>38.06</u>	7.93	<u>60.20</u>	<u>81.30</u>	<u>88.20</u>	<u>47.17</u>	<u>25.28</u>
	RDE	Best	75.94	90.14	94.12	67.56	51.44	67.68	82.47	87.36	40.06	<u>7.87</u>	65.35	83.95	89.90	50.88	28.08
20%	SSAN	Best	46.52	68.36	77.42	42.49	28.13	40.57	62.58	71.53	20.93	2.22	35.10	60.00	71.45	28.90	12.08
		Last	45.76	67.98	76.28	40.05	24.12	40.28	62.68	71.53	20.98	2.25	33.45	58.15	69.60	26.46	10.08
	IVT	Best	58.59	78.51	85.61	57.19	45.78	50.21	69.14	76.18	34.72	8.77	43.65	66.50	75.70	37.22	20.47
		Last	57.67	78.04	85.02	56.17	44.42	48.70	67.42	75.06	34.44	9.25	37.95	63.35	73.75	34.24	19.67
	IRRA	Best	69.74	87.09	92.20	62.28	45.84	60.76	78.26	84.01	35.87	6.80	58.75	81.90	88.25	46.38	24.78
		Last	69.44	87.09	92.04	62.16	45.70	60.58	78.14	84.20	35.92	6.91	54.00	77.15	88.55	43.20	22.53
	CLIP-C	Best	66.41	85.15	90.89	59.36	43.02	55.25	74.76	81.32	31.09	4.94	54.45	77.80	86.70	42.58	21.38
		Last	66.10	86.01	91.02	59.77	43.57	55.17	74.58	81.46	31.12	4.97	53.20	76.25	85.40	41.95	21.95
	DECL	Best	70.29	87.04	91.93	62.84	46.54	61.95	78.36	83.88	36.08	6.25	61.75	80.70	86.90	47.70	26.07
		Last	70.08	87.20	92.14	62.86	46.63	61.95	78.36	83.88	36.08	6.25	60.85	80.45	86.65	47.34	25.86
	RDE	Best	74.46	89.42	93.63	66.13	49.66	66.54	81.70	86.70	39.08	7.55	64.45	83.50	90.00	49.78	27.43
		Last	74.53	<u>89.23</u>	<u>93.55</u>	66.13	<u>49.63</u>	<u>66.51</u>	81.70	86.71	39.09	7.56	<u>63.85</u>	83.85	<u>89.45</u>	50.27	27.75
50%	SSAN	Best	13.43	31.74	41.89	14.12	6.91	18.83	37.70	47.43	9.83	1.01	19.40	39.25	50.95	15.95	6.13
		Last	11.31	28.07	37.90	10.57	3.46	17.06	37.18	47.85	6.58	0.39	14.10	33.95	46.55	11.88	4.04
	IVT	Best	50.49	71.82	79.81	48.85	36.60	43.03	61.48	69.56	28.86	6.11	39.70	63.80	73.95	34.35	18.56
		Last	42.02	65.04	73.72	40.49	27.89	36.57	54.83	62.91	24.30	5.08	28.55	52.05	62.70	26.82	13.97
	IRRA	Best	62.41	82.23	88.40	55.52	38.48	52.53	71.99	79.41	29.05	4.43	56.65	78.40	86.55	42.41	21.05
		Last	42.79	64.31	72.58	36.76	21.11	39.22	60.52	69.26	19.44	1.98	31.15	55.40	65.45	23.96	9.67
	CLIP-C	Best	64.02	83.66	89.38	57.33	40.90	51.60	71.89	79.31	28.76	4.33	53.45	76.80	85.50	41.43	21.17
		Last	63.97	83.74	89.54	57.35	40.88	51.49	71.99	79.32	28.77	4.37	52.35	76.35	85.25	40.64	20.45
	DECL	Best	65.22	83.72	89.28	57.94	41.39	<u>57.50</u>	75.09	<u>81.24</u>	<u>32.64</u>	<u>5.27</u>	<u>56.75</u>	<u>80.55</u>	<u>87.65</u>	<u>44.53</u>	23.61
		Last	65.09	83.58	89.26	57.89	41.35	57.49	75.10	81.23	32.63	5.26	55.00	80.50	86.50	43.81	23.31
	RDE	Best	71.33	87.41	91.81	63.50	47.36	63.76	79.53	84.91	37.38	6.80	62.85	83.20	89.15	47.67	23.97
		Last	<u>71.25</u>	<u>87.39</u>	<u>91.76</u>	63.59	47.50	63.76	79.53	84.91	37.38	6.80	62.85	83.20	89.15	47.67	<u>23.96</u>

Table 1. Performance comparison under different noise rates on three benchmarks. “Best” means choosing the best checkpoint on the validation set to test, and “Last” stands for choosing the checkpoint after the last training epoch to conduct inference. R-1,5,10 is an abbreviation for Rank-1,5,10 (%) accuracy. The best and second-best results are in **bold** and underline, respectively.

4.1.3 Implementation Details

As mentioned earlier, we adopt the pre-trained model CLIP [40] as our modality-specific encoders. In fairness, we use the same version of CLIP-ViT-B/16 as IRRA [24] to conduct experiments. During training, we introduce data augmentations to increase the diversity of the training data. Specifically, we utilize random horizontal flipping, random crop with padding, and random erasing to augment the training images. For training texts, we employ random masking, replacement, and removal for the word tokens as the data augmentation. Moreover, the input size of images is 384×128 and the maximum length of input word tokens is set to 77. We employ the Adam optimizer to train our model for 60 epochs with a cosine learning rate decay strategy. The initial learning rate is $1e - 5$ for the original model parameters of CLIP and the initial one for the network parameters of TSE is initialized to $1e - 3$. The batch size is 64. Following IRRA [24], we adopt an early training process with a gradually increasing learning rate. For hyperparameter settings, the margin value m of TAL is set to 0.1, the temperature parameter τ is set to 0.015, and the selection ratio \mathcal{R} is 0.3.

4.2. Comparison with State-of-the-Art Methods

In this section, we evaluate the performance of our RDE on three benchmarks under different scenarios. For a comprehensive comparison, we compare our method with several

state-of-the-art methods, including both ordinary methods and robust methods. Moreover, we use two synthetic noise levels (*i.e.*, noise rates), 20%, and 50%, to simulate the real-world scenario where the image-text pairs are not well-aligned. We randomly shuffle the text descriptions to inject NCs into the training data. We compare our RDE with five state-of-the-art baselines: SSAN [8], IVT [45], IRRA [24], DECL [36], and CLIP-C. SSAN, IVT, and IRRA are recent ordinary methods that are not designed for NC. DECL is a general framework that can enhance the robustness of image-text matching methods against NC. We use the model of IRRA as the base model of DECL for TIReID. CLIP-C is a strong baseline that fine-tunes the CLIP(ViT-B/16) model with only clean image-text pairs. We report the results of both the best checkpoint on the validation set and the last checkpoint to show the overfitting degree. Furthermore, we also evaluate our RDE on the original datasets without synthetic NC to demonstrate its superiority in Table 1. We compare our RDE with two local-matching methods: SSAN [8] and CFine [53]); and two global-matching methods: IVT [45] and IRRA [24]. More comparisons with other methods are provided in the supplementary material.

From Table 1, one can see that our RDE achieves state-of-the-art performance on three datasets and we can draw three observations: (1) On the datasets with synthetic NC, the ordinary methods suffer from remarkable performance degradation or poor performance as the noise rate increases. In contrast, our RDE achieves the best results on all met-

rics. Moreover, by comparing the ‘Best’ performance with the ‘Last’ ones in Table 1, we can see that our RDE can effectively prevent the performance deterioration caused by overfitting against NC. (2) Compared with the robust framework DECL and the strong baseline CLIP-C, our RDE also shows obvious advantages, which indicates that our solution against NC is effective and superior in TIReID. For instance, on CUHK-PEDES under 50% noise, our RDE achieves 71.33%, 87.41%, and 91.81% in terms of Rank-1, 5, 10 on the ‘Best’ rows, respectively, which surpasses the best baseline DECL by a large margin, *i.e.*, +6.11%, +3.69%, and +2.53%, respectively. (3) On the datasets without synthetic NC, our RDE outperforms all baselines by a large margin. Specifically, RDE achieves performance gains of +2.56%, +4.22%, and +5.15% in terms of Rank-1 compared with the best baseline IRRA on three datasets, respectively, demonstrating the effectiveness and advantages of our method.

4.3. Ablation Study

In this section, we conduct ablation studies on the CUHK-PEDES dataset with 50% noise to investigate the effects and contributions of each proposed component in RDE. We compare different combinations of our components in Table 2. From the experimental results, we could draw the following observation: (1) RDE achieves the best performance by using both BGE and TSE for joint inference, which demonstrates that these two modules are complementary and effective. (2) RDE benefits from CCD, which can enhance the robustness and alleviate the overfitting effect caused by NC. (3) Our TAL outperforms the widely-used Triplet Ranking Loss (TRL) and SDM loss [24], which demonstrates the superior stability and robustness of our TAL against NC.

No.	S^b	S^e	CCD	Loss	R-1	R-5	R-10	mAP	mINP
#1	✓	✓	✓	TAL	71.33	87.41	91.81	63.50	47.36
#2	✓	✓	✓	TRL	6.40	16.08	22.14	6.53	2.51
#3	✓	✓	✓	TRL-S	67.38	85.35	90.64	60.04	43.60
#4	✓	✓	✓	SDM	69.33	86.99	91.68	61.99	45.34
#5		✓	✓	TAL	70.70	86.60	91.16	62.67	46.19
#6	✓		✓	TAL	69.07	86.09	91.13	61.69	45.40
#7	✓	✓		TAL	63.11	81.04	87.22	55.42	38.68

Table 2. Ablation studies on the CHUK-PEDES dataset.

4.4. Parametric Analysis

To study the impact of different hyperparameter settings on performance, we perform sensitivity analyses for two key hyperparameters (*i.e.*, m and τ) on the CHUK-PEDES dataset with 50% noise. From Figure 4, we can see that: (1) Too large or too small m will lead to suboptimal performance. We choose $m = 0.1$ in all our experiments. (2) Too small τ will cause training failure, while the increasing τ will gradually decrease the separability (hardness) of

positive and negative pairs for suboptimal performance. We choose $\tau = 0.015$ in all our experiments.

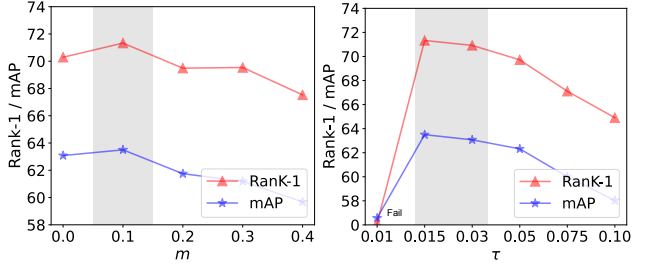


Figure 4. Variation of performance with different m and τ .

4.5. Robustness Study

In this section, we provide some visualization results during cross-modal training to verify the robustness and effectiveness of our method. As shown in Figure 5, one can clearly see that our RDE not only achieves excellent performance under noise but also effectively alleviates noise overfitting.

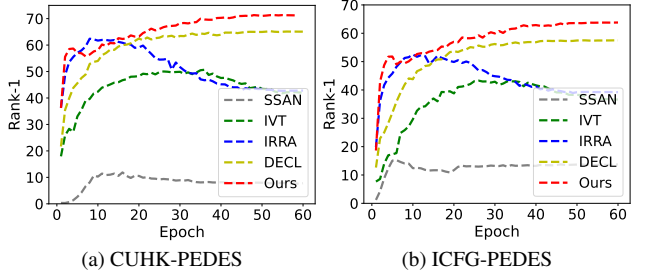


Figure 5. Test performance (Rank-1) versus epochs on the CHUK-PEDES and ICFG-PEDES datasets with 50% noise.

5. Conclusion

In this paper, we reveal and study a novel challenging problem of noisy correspondence (NC) in TIReID, which violates the common assumption of existing methods that image-text data is perfectly aligned. To this end, we propose a robust method, *i.e.*, RDE, to effectively handle the revealed NC problem and achieve superior performance. Extensive experiments are conducted on three datasets to comprehensively demonstrate the superiority and robustness of RDE both with and without synthetic NCs.

Acknowledgments

This work was supported in part by NSFC under Grant U21B2040, 62176171, 62372315, and 62102274, in part by Sichuan Science and Technology Program under Grant 2022YFH0021 and 2023ZYD0143; in part by Chengdu Science and Technology Project under Grant 2023-XT00-00004-GX; in part by the SCU-LuZhou Sciences and Technology Cooperation Program under Grant 2023CDLZ-16; in part by the Fundamental Research Funds for the Central Universities under Grant CJ202303 and YJ202140.

Supplementary Material: Noisy-Correspondence Learning for Text-Image Person Re-identification

In this supplementary material, we provide additional information for RDE. More specifically, we first give detailed proof and derivation for lemmas and gradients in Appendix A. In Appendix B, we detail the used datasets and the compared baselines. In Appendix C, to further verify the robustness of RDE, we provide the re-identification performance on three benchmark datasets under extremely high noise rate, *i.e.*, 80%. Besides, in Appendix D, we provide more comparison results compared with state-of-the-art methods to comprehensively verify the superiority of our RDE. In Appendix E, we explore the impact of different selection ratios (\mathcal{R}) on performance. In Appendix F, we provide a more ablation analysis. In Appendix G, we provide a large number of real noisy examples existing in the three public datasets to conduct a case study, thus emphasizing our motivation. We also provide a more comprehensive robustness analysis to verify the robustness of RDE in Appendix H. Finally, in Appendix I, we provide some qualitative results to illustrate the advantages of our RDE.

A. Proof and Derivation

A.1. Proof for Lemma 1 (Lemma 2)

Given an input image-text pair (I_i, T_i) in a mini-batch \mathbf{x} , TAL is defined as:

$$\begin{aligned} \mathcal{L}_{tal}(I_i, T_i) = & [m - S_{i2t}^+(I_i) + \tau \log(\sum_{j=1}^K q_{ij} \exp(S(I_i, T_j)/\tau))]_+ \\ & + [m - S_{t2i}^+(T_i) + \tau \log(\sum_{j=1}^K q_{ji} \exp(S(I_j, T_i)/\tau))]_+, \end{aligned} \quad (11)$$

where m is a positive margin coefficient, τ is a temperature coefficient to control hardness, $S(I_i, T_j) \in \{S_{ij}^b, S_{ij}^t\}$, $[x]_+ \equiv \max(x, 0)$, $\exp(x) \equiv e^x$, $q_{ij} = 1 - l_{ij}$, and K is the size of \mathbf{x} . From Lemma 1, as $\tau \rightarrow 0$, TAL is close to TRL and focuses more on hard negatives. Since multiple positive pairs from the same identity may appear in the mini-batch, $S_{i2t}^+(I_i) = \sum_{j=1}^K \alpha_{ij} S(I_i, T_j)$ is the weighted average similarity of positive pairs for image I_i , where $\alpha_{ij} = \frac{l_{ij} \exp(S(I_i, T_j)/\tau)}{\sum_{k=1}^N l_{ik} \exp(S(I_i, T_k)/\tau)}$. And, $S_{i2t}^+(T_i)$ is similar to the definition of $S_{i2t}^+(I_i)$.

Lemma 2 TAL is the upper bound of TRL, *i.e.*,

$$\begin{aligned} \mathcal{L}_{trl}(I_i, T_i) = & [m - S_{i2t}^+(I_i) + S(I_i, \hat{T}_i)]_+ \\ & + [m - S_{t2i}^+(T_i) + S(\hat{I}_i, T_i)]_+ \leq \mathcal{L}_{tal}(I_i, T_i), \end{aligned} \quad (12)$$

where $\hat{T}_i \in \mathbf{T}_i = \{T_j | l_{ij} = 0, \forall j \in \{1, 2, \dots, K\}\}$ is the hardest negative text for image I_i and $\hat{I}_i \in \mathbf{I}_i = \{I_j | l_{ji} = 0, \forall j \in \{1, 2, \dots, K\}\}$ is the hardest negative image for text I_i , respectively.

Proof 1 To prove Equation (12), we first take the image-to-text direction as an example. For $S(I_i, \hat{T}_i)$ in Equation (12), we have that

$$\begin{aligned} S(I_i, \hat{T}_i) = & \max_{T_j \in \mathbf{T}_i} (S(I_i, T_j)) \\ = & \max_{T_j \in \mathbf{T}_i} \left(\tau \log \exp(S(I_i, T_j))^{\frac{1}{\tau}} \right) \\ = & \tau \log \left(\max_{T_j \in \mathbf{T}_i} \left(\exp(S(I_i, T_j))^{\frac{1}{\tau}} \right) \right) \\ \leq & \tau \log \left(\sum_{T_j \in \mathbf{T}_i} \exp(S(I_i, T_j)/\tau) \right) \\ \leq & \tau \log \left(\sum_{j=1}^K q_{ij} \exp(S(I_i, T_j)/\tau) \right), \end{aligned} \quad (13)$$

where $q_{ij} = 1 - l_{ij}$. Based on Equation (13), we have that

$$\begin{aligned} & [m - S_{i2t}^+(I_i) + \tau \log(\sum_{j=1}^K q_{ij} \exp(S(I_i, T_j)/\tau))]_+ \\ & \geq [m - S_{i2t}^+(I_i) + S(I_i, \hat{T}_i)]_+. \end{aligned} \quad (14)$$

Similarly, in the text-to-image direction, we have that

$$\begin{aligned} & [m - S_{t2i}^+(T_i) + \tau \log(\sum_{j=1}^K q_{ji} \exp(S(I_j, T_i)/\tau))]_+ \\ & \geq [m - S_{t2i}^+(T_i) + S(\hat{I}_i, T_i)]_+. \end{aligned} \quad (15)$$

Thus, combining Equation (14) and Equation (15), we can get $\mathcal{L}_{trl}(I_i, T_i) \leq \mathcal{L}_{tal}(I_i, T_i)$. This completes the proof.

A.2. Derivation for Gradient

In this appendix, we provide more details of gradient derivation. For ease of representation and analysis, we only consider one direction like [30] since image-to-text retrieval and text-to-image retrieval are symmetrical. Besides, we suppose that there is only one paired text for each image in the mini-batch. Thus, TRL, TRL-S, and TAL are simplified as follows:

$$\begin{aligned}\mathcal{L}_{trl}(I_i, T_i) &= [m - \mathbf{v}_i^\top \hat{\mathbf{t}}_i + \mathbf{v}_i^\top \hat{\mathbf{t}}_i]_+, \\ \mathcal{L}_{trls}(I_i, T_i) &= \sum_{j \neq i}^K [m - \mathbf{v}_i^\top \hat{\mathbf{t}}_i + \mathbf{v}_i^\top \hat{\mathbf{t}}_j]_+, \\ \mathcal{L}_{tal}(I_i, T_i) &= \left[m - \mathbf{v}_i^\top \hat{\mathbf{t}}_i + \tau \log \left(\sum_{j \neq i}^K e^{(\mathbf{v}_i^\top \hat{\mathbf{t}}_j / \tau)} \right) \right]_+, \end{aligned} \quad (16)$$

where $\hat{\mathbf{t}}_i$, $\hat{\mathbf{t}}_j$ and \mathbf{t}_i are the hardest negative sample, negative sample, and positive sample of the anchor sample \mathbf{v}_i , respectively. These ℓ_2 -normalized features are embedded by the modality-specific models, *i.e.*, $f_{\Theta_v}(\cdot)$ and $f_{\Theta_t}(\cdot)$. Due to the truncation operation $[x]_+$, we only discuss the case of $\mathcal{L} > 0$ that could generate gradients. For TRL, the gradients to the parameters Θ_v and Θ_t are:

$$\begin{aligned}\frac{\partial \mathcal{L}_{trl}}{\partial \Theta_v} &= \frac{\partial \mathcal{L}_{trl}}{\partial \mathbf{v}_i} \frac{\partial \mathbf{v}_i}{\partial \Theta_v}, \\ \frac{\partial \mathcal{L}_{trl}}{\partial \Theta_t} &= \frac{\partial \mathcal{L}_{trl}}{\partial \hat{\mathbf{t}}_i} \frac{\partial \hat{\mathbf{t}}_i}{\partial \Theta_t} + \frac{\partial \mathcal{L}_{trl}}{\partial \mathbf{t}_i} \frac{\partial \mathbf{t}_i}{\partial \Theta_t}. \end{aligned} \quad (17)$$

Since the learning of normalized features can be viewed as the movement process of points on a unit hyperplane, we only consider the loss gradients with respect to \mathbf{v}_i , $\hat{\mathbf{v}}_i$, and \mathbf{t}_i are:

$$\frac{\partial \mathcal{L}_{trl}}{\partial \mathbf{v}_i} = \hat{\mathbf{t}}_i - \mathbf{t}_i, \quad \frac{\partial \mathcal{L}_{trl}}{\partial \hat{\mathbf{t}}_i} = -\mathbf{v}_i, \quad \frac{\partial \mathcal{L}_{trl}}{\partial \mathbf{t}_i} = \mathbf{v}_i. \quad (18)$$

For TRL-S, the gradients to the parameters Θ_v and Θ_t are:

$$\begin{aligned}\frac{\partial \mathcal{L}_{trls}}{\partial \Theta_v} &= \frac{\partial \mathcal{L}_{trls}}{\partial \mathbf{v}_i} \frac{\partial \mathbf{v}_i}{\partial \Theta_v}, \\ \frac{\partial \mathcal{L}_{trls}}{\partial \Theta_t} &= \sum_{j \in \mathcal{Z}} \frac{\partial \mathcal{L}_{trls}}{\partial \hat{\mathbf{t}}_j} \frac{\partial \hat{\mathbf{t}}_j}{\partial \Theta_t} + \frac{\partial \mathcal{L}_{trls}}{\partial \mathbf{t}_i} \frac{\partial \mathbf{t}_i}{\partial \Theta_t}. \end{aligned} \quad (19)$$

Thus, for \mathbf{v}_i , \mathbf{v}_j , and \mathbf{t}_i , the gradients are:

$$\begin{aligned}\frac{\partial \mathcal{L}_{trls}}{\partial \mathbf{v}_i} &= \sum_{j \in \mathcal{Z}} (\mathbf{t}_j - \mathbf{t}_i), \quad \frac{\partial \mathcal{L}_{trls}}{\partial \mathbf{t}_j} = \mathbf{v}_i, \forall j \in \mathcal{Z}, \\ \frac{\partial \mathcal{L}_{trls}}{\partial \mathbf{t}_i} &= -\sum_{j \in \mathcal{Z}} \mathbf{v}_i = -|\mathcal{Z}| \mathbf{v}_i, \end{aligned} \quad (20)$$

where $\mathcal{Z} = \{z \mid [m - S(I_i, T_i) + S(I_i, T_z)]_+ > 0, z \neq i, z \in \{0, \dots, K\}\}$. For our TAL, the gradients to the pa-

rameters Θ_v and Θ_t are:

$$\begin{aligned}\frac{\partial \mathcal{L}_{tal}}{\partial \Theta_v} &= \frac{\partial \mathcal{L}_{tal}}{\partial \mathbf{v}_i} \frac{\partial \mathbf{v}_i}{\partial \Theta_v}, \\ \frac{\partial \mathcal{L}_{tal}}{\partial \Theta_t} &= \sum_{j \neq i} \frac{\partial \mathcal{L}_{tal}}{\partial \hat{\mathbf{t}}_j} \frac{\partial \hat{\mathbf{t}}_j}{\partial \Theta_t} + \frac{\partial \mathcal{L}_{tal}}{\partial \mathbf{t}_i} \frac{\partial \mathbf{t}_i}{\partial \Theta_t}. \end{aligned} \quad (21)$$

Thus, the gradients for \mathbf{v}_i , \mathbf{v}_j \mathbf{t}_i are:

$$\begin{aligned}\frac{\partial \mathcal{L}_{tal}}{\partial \mathbf{v}_i} &= \sum_{j \neq i}^K \beta_j \mathbf{t}_j - \mathbf{t}_i = \sum_{j \neq i}^K \beta_j (\mathbf{t}_j - \mathbf{t}_i), \\ \frac{\partial \mathcal{L}_{tal}}{\partial \mathbf{t}_i} &= -\mathbf{v}_i, \quad \frac{\partial \mathcal{L}_{tal}}{\partial \mathbf{t}_j} = \beta_j \mathbf{v}_i, \end{aligned} \quad (22)$$

where $\beta_j = \frac{\exp(\mathbf{v}_i^\top \mathbf{t}_j / \tau)}{\sum_{k \neq i}^K \exp(\mathbf{v}_i^\top \mathbf{t}_k / \tau)}$.

B. Dataset and Baseline Description

B.1. Datasets.

To verify the effectiveness and superiority of RDE, we use three widely-used image-text person datasets to conduct experiments. A brief introduction of these datasets is given as follows:

- **CHUK-PEDES** [27] is the first large-scale benchmark to dedicate TIReID, which includes 40,206 person images and 80,412 text descriptions for 13,003 unique identities. We follow the official data split to conduct experiments, *i.e.*, 11,003 identities for training, 1,000 identities for validation, and the rest of the 1,000 identities for testing.
- **ICFG-PEDES** [8] is a widely-used benchmark collected from the MSMT17 dataset [51] and consists of 54,522 images for 4,102 unique persons and each image has a corresponding textual description. We follow the data split used by most TIReID methods [24, 45], *i.e.*, a training set with 3,102 identifies and a validation set with 1,000 identities. Note that we uniformly used the validation performance as the test performance due to its lack of a test set.
- **RSTPReid** [62] is another benchmark dataset constructed from the MSMT17 dataset [51] for TIReID. RSTPReid contains 20,505 images for 4,101 identities, wherein each person has 5 images and each image is paired with 2 text descriptions. Following the official data split, we use 3,701 identities for training, 200 identities for validation, and the remaining 200 identities for testing.

B.2. Baselines.

To verify the effectiveness and robustness of our method in the NC scenario, we provide the comparison results with 5 baselines that have published code. We introduce each baseline as follows:

- **SSAN**¹ [8] is a local-matching method for TIReID, which mainly benefits from a proposed multiview non-local net-

¹<https://github.com/zifyloo/SSAN>

Noise	Methods		CUHK-PEDES					ICFG-PEDES					RSTPReid				
			R-1	R-5	R-10	mAP	mINP	R-1	R-5	R-10	mAP	mINP	R-1	R-5	R-10	mAP	mINP
80%	SSAN	Best	0.18	0.83	1.45	0.47	0.24	0.28	0.99	1.90	0.27	0.15	0.65	3.25	5.95	1.30	0.70
		Last	0.13	0.67	1.46	0.42	0.21	0.18	1.01	1.77	0.25	0.14	0.65	2.95	5.85	1.32	0.68
	IVT	Best	34.03	55.49	66.16	33.90	23.29	21.10	37.10	45.64	13.68	2.32	15.15	30.00	40.50	14.98	7.79
		Last	10.61	23.81	31.38	11.13	5.72	5.64	12.48	17.15	4.00	0.69	4.95	13.55	19.75	6.07	2.85
	IRRA	Best	38.63	56.69	64.18	34.60	21.84	28.19	44.14	51.27	14.36	1.41	29.65	46.65	54.50	23.77	11.32
		Last	9.06	19.69	25.65	8.26	3.18	8.68	18.76	24.50	3.65	0.27	8.15	21.00	29.05	7.28	2.40
	CLIP-C	Best	57.38	78.05	84.97	51.08	34.83	44.84	65.24	73.27	24.27	3.42	47.80	72.70	81.75	37.50	18.09
		Last	57.05	78.09	85.07	51.14	<u>35.05</u>	44.65	65.26	73.45	24.20	3.44	44.60	70.75	80.20	35.67	17.09
	DECL	Best	47.90	71.57	80.17	44.51	29.86	40.53	61.49	69.84	21.78	2.97	48.15	72.20	80.75	37.31	18.83
		Last	46.57	70.19	78.48	42.93	27.91	39.91	61.16	69.51	21.56	2.89	45.85	71.05	81.00	35.34	16.35
	RDE	Best	64.99	<u>83.15</u>	<u>89.52</u>	57.84	41.07	56.02	<u>74.00</u>	80.62	<u>30.67</u>	<u>4.60</u>	53.40	<u>76.70</u>	<u>85.55</u>	<u>39.71</u>	<u>18.28</u>
		Last	<u>64.91</u>	83.20	89.54	<u>57.83</u>	41.07	<u>55.96</u>	74.09	<u>80.61</u>	30.79	4.62	<u>52.35</u>	76.85	84.90	39.92	17.72

Table 3. Performance comparison under 80% noise rate on three benchmarks. “Best” means choosing the best checkpoint on the validation set to test, and “Last” stands for choosing the checkpoint after the last training epoch to conduct inference. R-1,5,10 is an abbreviation for Rank-1,5,10 (%) accuracy. The best and second-best results are in **bold** and underline, respectively.

work that could capture the local relationships, thus establishing better correspondences between body parts and noun phrases. Besides, SSAN also exploits a compound ranking loss to make an effective reduction of the intra-class variance in textual features.

- **IVT**² [45] is an implicit visual-textual framework, which belongs to the global-matching method. To explore fine-grained alignments, IVT utilizes two implicit semantic alignment paradigms, *i.e.*, multi-level alignment (MLA) and bidirectional mask modeling (BMM). MLA aims to see “finer” by exploring local and global alignments from three-level matchings. BMM aims to see “more” by mining more semantic alignments from random masking for both modalities.
- **IRRA**³ [24] is a recent state-of-art global-matching method that could learn relations between local visual-textual tokens and enhances global alignments without requiring additional prior supervision. IRRA exploits a novel similarity distribution matching to minimize the KL divergence between the similarity distributions and the normalized label matching distributions for better performance.
- **CLIP-C** is a quite strong baseline that fine-tunes the original CLIP⁴ model with only clean image-text pairs. We use the same version as IRRA, *i.e.*, ViTB/16, for a fair comparison and use InfoNCE loss [35] to optimize the model.
- **DECL**⁵ [36] is an effective robust image-text matching framework, which utilizes the cross-modal evidential learning paradigm to capture and leverage the uncertainty brought by noise to isolate the noisy pairs. Since TIReID can be treated as the sub-task of instance-level image-text

matching, DECL also can be used to ease the negative impact of NCs in TIReID. In this paper, we exploit the used model of IRRA [24] as the base model of DECL for robust TIReID.

C. The Results under Extreme Noise

To further verify the effectiveness and robustness of our method, we report comparison results under extremely high noise, *i.e.*, 80%. From the results in Table 3, one can see that our RDE achieves the best performance and can effectively alleviate the performance degradation caused by noise overfitting. For example, compared with the ‘Best’ rows, our RDE surpasses the best baselines by +7.56%, +5.95%, and +3.5% in terms of Rank-1 on the three datasets, respectively.

D. More Comparisons

In this section, we follow the organization of IRRA [24] and provide more comparative experimental results on three benchmarks in Tables 4 to 6. From the results, our RDE achieves the best results and exceeds the best baselines, *i.e.*, +0.92%, +2.63%, and +0.15% in terms of Rank-1 on three datasets, respectively.

E. Study on the Selection Ratio

Figure 6 shows the variation of performance with different selection ratio \mathcal{R} . From the figure, one can see that too large or too small \mathcal{R} will cause suboptimal performance. We think that a small \mathcal{R} will cause too much information loss and poor embedding presentations, while too large will focus on too many meaningless features. For this reason, we recommend \mathcal{R} to be set between 0.3~0.5. Thus, \mathcal{R} is set to 0.3 in all our experiments.

²<https://github.com/TencentYoutuResearch/PersonRetrieval-IVT>

³<https://github.com/anosorae/IRRA>

⁴<https://github.com/openai/CLIP>

⁵<https://github.com/QinYang79/DECL>

Methods	Ref.	Image Enc.	Text Enc.	R-1	R-5	R-10	mAP	mINP
CMPM/C [60]	ECCV'18	RN50	LSTM	49.37	-	79.27	-	-
TIMAM [41]	ICCV'19	RN101	BERT	54.51	77.56	79.27	-	-
ViTAA [48]	ECCV'20	RN50	LSTM	54.92	75.18	82.90	51.60	-
NAFS [16]	arXiv'21	RN50	BERT	59.36	79.13	86.00	54.07	-
DSSL [62]	MM'21	RN50	BERT	59.98	80.41	87.56	-	-
SSAN [8]	arXiv'21	RN50	LSTM	61.37	80.15	86.73	-	-
LapScore [52]	ICCV'21	RN50	BERT	63.40	-	87.80	-	-
ISANet [54]	arXiv'22	RN50	LSTM	63.92	82.15	87.69	-	-
LBUL [50]	MM'22	RN50	BERT	64.04	82.66	87.22	-	-
Han et al. 2021	BMVC'21	CLIP-RN101	CLIP-Xformer	64.08	81.73	88.19	60.08	-
SAF [28]	ICASSP'22	ViT-Base	BERT	64.13	82.62	88.40	-	-
TIPCB [5]	Neuro'22	RN50	BERT	64.26	83.19	89.10	-	-
CAIBC [49]	MM'22	RN50	BERT	64.43	82.87	88.37	-	-
AXM-Net [13]	MM'22	RN50	BERT	64.44	80.52	86.77	58.73	-
LGUR [42]	MM'22	DeiT-Small	BERT	65.25	83.12	89.00	-	-
IVT [45]	ECCVW'22	ViT-Base	BERT	65.59	83.11	89.21	-	-
CFine [53]	TIP'23	CLIP-ViT	BERT	69.57	85.93	91.15	-	-
IRRA [24]	CVPR'23	CLIP-ViT	CLIP-Xformer	73.38	89.93	93.71	66.13	50.24
BiLMa [15]	ICCVW'23	CLIP-ViT	CLIP-Xformer	74.03	89.59	93.62	66.57	-
PBSL [44]	ACMMM'23	RN50	BERT	65.32	83.81	89.26	-	-
BEAT[33]	ACMMM'23	RN101	BERT	65.61	83.45	89.54	-	-
LCR ² S [55]	ACMMM'23	RN50	TextCNN	67.36	84.19	89.62	59.24	-
DCEL [29]	ACMMM'23	CLIP-ViT	CLIP-Xformer	75.02	90.89	94.52	-	-
UniPT [43]	ICCV'23	CLIP-ViT	CLIP-Xformer	68.50	84.67	-	-	-
RaSa [2]	IJCAI'23	ALBEFF	ALBEFF	76.51	90.29	94.25	69.38	-
RaSa _{TCL} [2]	IJCAI'23	TCL	TCL	73.23	89.20	93.32	66.43	-
TBPS [4]	Arxiv'23	CLIP-ViT	CLIP-Xformer	73.54	88.19	92.35	65.38	-
Our RDE	-	CLIP-ViT	CLIP-Xformer	75.94	90.14	94.12	67.56	51.44

Table 4. Performance comparisons on the CUHK-PEDES dataset. The best results are in **bold**.

Methods	R-1	R-5	R-10	mAP	mINP
Dual Path [61]	38.99	59.44	68.41	-	-
CMPM/C [60]	43.51	65.44	74.26	-	-
ViTAA [48]	50.98	68.79	75.78	-	-
SSAN [8]	54.23	72.63	79.53	-	-
IVT [45]	56.04	73.60	80.22	-	-
ISANet [54]	57.73	75.42	81.72	-	-
CFine [53]	60.83	76.55	82.42	-	-
IRRA [24]	63.46	80.25	85.82	38.06	7.93
BiLMa [15]	63.83	80.15	85.74	38.26	-
PBSL [44]	57.84	75.46	82.15	-	-
BEAT[33]	58.25	75.92	81.96	-	-
LCR ² S [55]	57.93	76.08	82.40	38.21	-
DCEL [29]	64.88	81.34	86.72	-	-
UniPT [43]	60.09	76.19	-	-	-
RaSa [2]	65.28	80.40	85.12	41.29	-
RaSa _{TCL} [2]	63.29	79.36	84.36	39.23	-
TBPS [4]	65.05	80.34	85.47	39.83	-
Our RDE	67.68	82.47	87.36	40.06	7.87

Table 5. Performance comparisons on the ICFG-PEDES dataset. The best results are in **bold**. ‘*’ indicates our reproducible results.

Methods	R-1	R-5	R-10	mAP	mINP
DSSL [62]	39.05	62.60	73.95	-	-
SSAN [8]	43.50	67.80	77.15	-	-
LBUL [50]	45.55	68.20	77.85	-	-
IVT [45]	46.70	70.00	78.80	-	-
CFine [53]	50.55	72.50	81.60	-	-
IRRA [24]	60.20	81.30	88.20	47.17	25.28
BiLMa [15]	61.20	81.50	88.80	48.51	-
PBSL [44]	47.80	71.40	79.90	-	-
BEAT[33]	48.10	73.10	81.30	-	-
LCR ² S [55]	54.95	76.65	84.70	40.92	-
DCEL [29]	61.35	83.95	90.45	-	-
RaSa [2]	66.90	86.50	91.35	52.31	-
RaSa _{TCL} [2]	65.20	84.05	89.85	50.14	-
TBPS [4]	61.95	83.55	88.75	48.26	-
Our RDE	65.35	83.95	89.90	50.88	28.08

Table 6. Performance comparisons on the RSTPReid dataset. The best results are in **bold**. ‘*’ indicates our reproducible results.

F. Ablation Study

F.1. Ablation analysis for TSE

To verify the design rationality of TSE in our RDE, we conduct dedicated ablation experiments on TSE. The results are

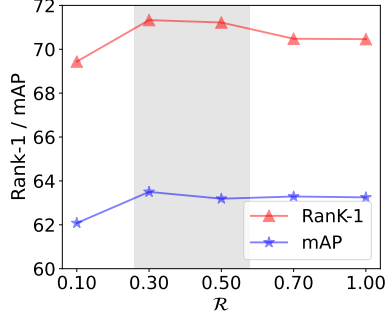


Figure 6. Variation of performance with different $\mathcal{R} \in [0, 1]$.

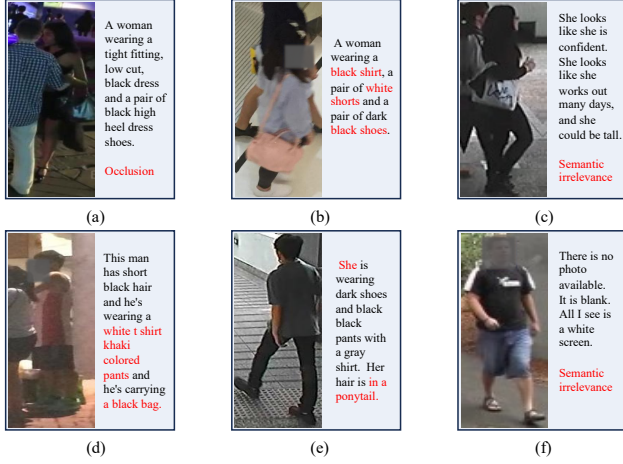


Figure 7. The examples of noisy correspondence identified by CCD on the CUHK-PEDES dataset.



Figure 8. The examples of noisy correspondence identified by CCD on the ICFG-PEDES dataset.

reported in Table 7. In the table, TSE' means that the token features encoded by CLIP are directly used for aggregation to obtain the embedding representations instead of conducting embedding transformation. Also, we show the impact of different pooling strategies on performance. From the

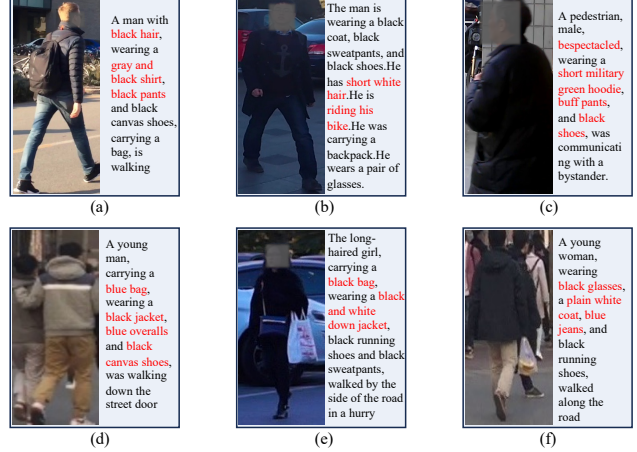


Figure 9. The examples of noisy correspondence identified by CCD on the RSTPReid dataset.

results, our standard version of TSE obtains the best performance, *i.e.*, conducting the embedding transformation and using the max-pooling strategy to obtain the TSE representations.

Methods	Pool	R-1	R-5	R-10	mAP	mINP
TSE'	Avg.	67.22	84.96	90.03	60.22	43.84
TSE'	TopK.	67.35	85.36	90.51	60.21	43.54
TSE'	Max.	67.46	85.17	90.58	60.11	43.45
TSE	Avg.	67.43	85.19	90.50	60.42	43.97
TSE	TopK.	68.27	86.03	90.79	60.95	44.37
TSE	Max.	71.33	87.41	91.81	63.50	47.36

Table 7. Performance comparisons with state-of-the-art methods on the RSTPReid dataset. 'Avg.', 'TopK.', and 'Max.' indicate the use of average-pooling, topK-pooling (K=10), and max-pooling strategies, respectively.

Noise	No.	S^b	S^t	CCD	Loss	R-1	R-5	R-10	mAP	mINP
80%	#1	✓	✓	✓	TAL	64.99	83.15	89.52	57.84	41.07
	#2	✓	✓	✓	TRL	2.18	6.45	10.48	2.65	0.83
	#3	✓	✓	✓	TRL-S	51.62	74.53	82.21	46.15	30.12
	#4	✓	✓	✓	SDM	58.32	79.03	85.79	51.27	34.00
	#5	✓	✓	✓	TAL	63.56	82.59	88.84	56.69	39.71
	#6	✓	✓	✓	TAL	61.70	81.61	87.95	55.11	38.34
	#7	✓	✓	✓	TAL	41.03	62.62	71.99	37.29	23.54

Table 8. Ablation studies on the CHUK-PEDES dataset.

F.2. Ablation study on High Noise

In this appendix, we provide more ablation studies on the CUHK-PEDES dataset to investigate the effects and contributions of each proposed component in RDE. The experimental results are shown in Table 8. The observations

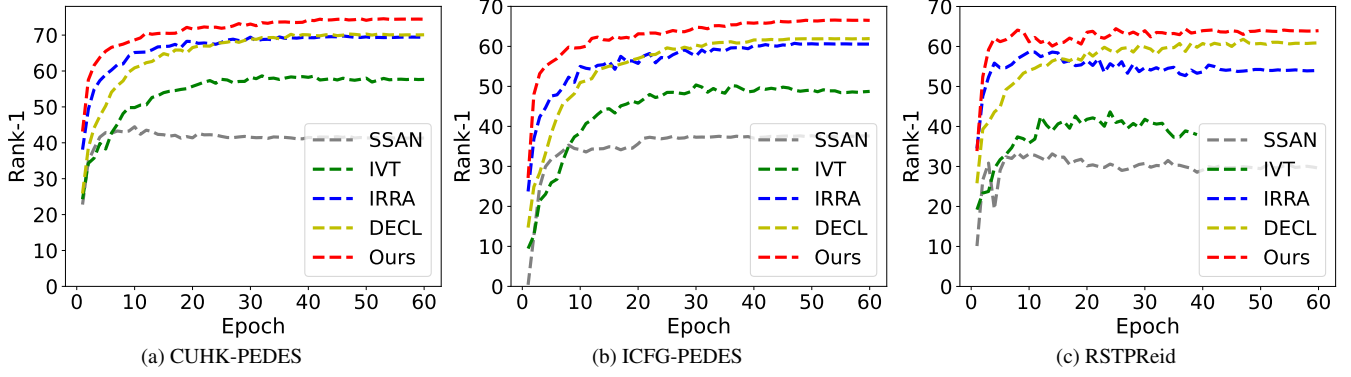


Figure 10. Test performance (Rank-1) versus epochs on three datasets with 20% noise.

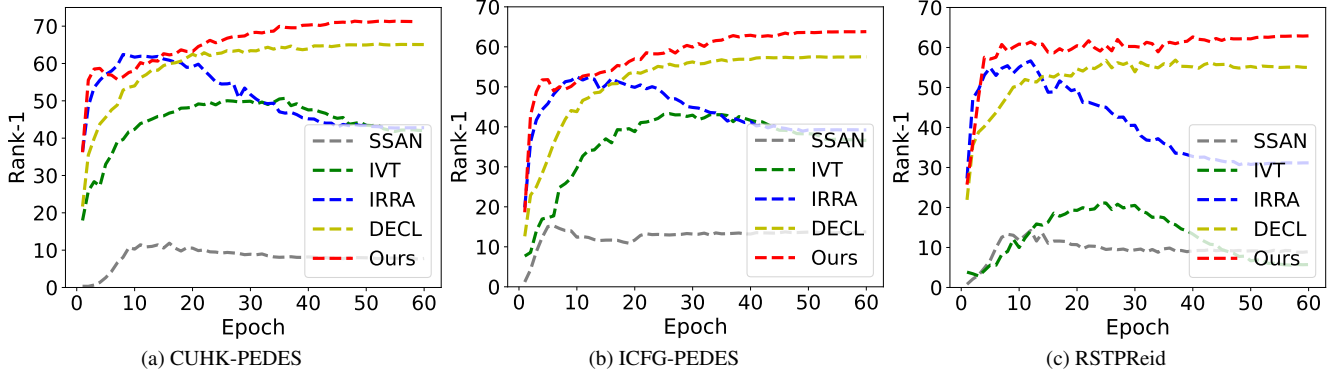


Figure 11. Test performance (Rank-1) versus epochs on three datasets with 50% noise.

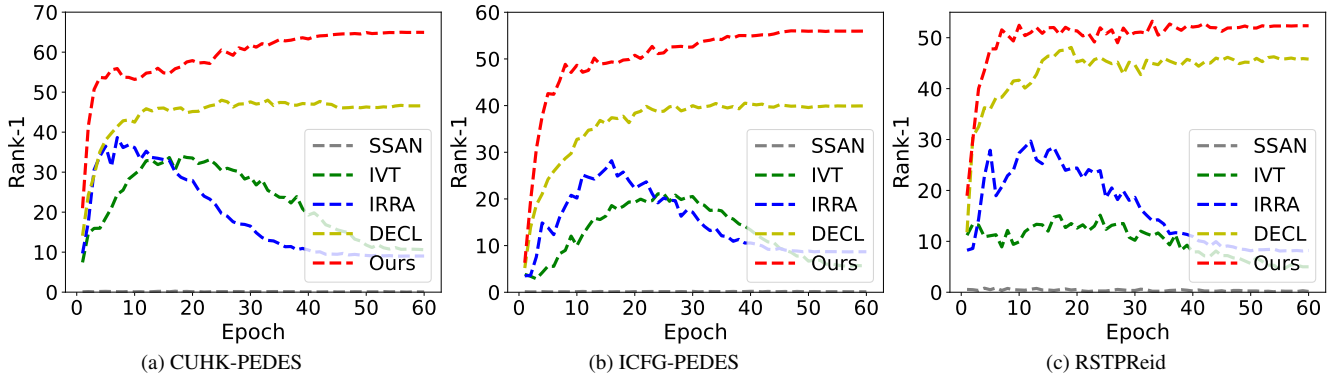


Figure 12. Test performance (Rank-1) versus epochs on three datasets with 80% noise.

and conclusions are consistent with those in the main text, which also demonstrate the effectiveness of our method.

G. Case Study

In this section, we show a large number of real examples of noisy pairs in three public datasets without synthetic NCs in Figures 1 and 2 and ??, which are identified by CCD. Note that for privacy and security, the face areas of people in all images are **blurred**. From these examples,

one can see that there are various reasons for noisy correspondences, *e.g.*, occlusion (*e.g.*, Figure 7(a,b)), lighting (*e.g.*, Figure 8(f)), and inaccurate noisy text descriptions (*e.g.*, Figure 7(c,f) and Figure 9(a-f)). But all in all, these noisy pairs are real in these datasets and actually **break the implicit assumption** that all training image-text pairs are aligned correctly and perfectly at an instance level. Thus, we reveal the noisy correspondence problem in TIReID and propose a robust method, *i.e.*, RDE, to particularly address it.

(a) A woman walking visible from the back is wearing a white shirt, black pants and has a green bag slung over her back and carrying a black object in her right hand.



(b) The pedestrian with long, dark hair carries a backpack. She wears a loose top, denim bottoms, and sandals.



(c) This person wearing the sneakers and dark hoodie is walking with a large shoulder bag.



(d) This person has a white band in their hair he or she is wearing a pancho in salmon color with a yellow bend on the bottom as well as a dark tight pants and dark shoes.



Figure 13. Comparison of top-10 retrieved results on the CUHK-PEDES dataset between the baseline IRRA (the first row) and our RDE (the second row) for each text query. The matched and mismatched person images are marked with red and blue rectangles, respectively. All face areas of people in images are **blurred** for privacy and security.

H. Robustness Study

For a comprehensive robustness analysis, we provide more performance curves versus epochs in Figures 10 to 12. It can be seen from the Figure 10 that when the noise rate is 20%, each baseline shows a certain degree of robustness, and there is no obvious performance degradation due to over-fitting noisy pairs. However, as the noise rate increases, the non-robust methods (SSAN, IVT, and IRRA) all show a curve that first rises and then falls. This tendency is caused by the memorization effect that DNNs tend to learn simple patterns before fitting noisy samples. Besides, we can also find that when the noise rate is 80%, SSAN fails and other non-robust methods (IVT and IRRA) also have a serious performance drop. By contrast, thanks to the CCD and TAL, our RDE can learn accurate visual-semantic associations by obtaining confident clean training image-text pairs, which can effectively and directly prevent over-fitting noisy pairs, thus achieving robust cross-modal learning. From these figures, our method not only exhibits strong robustness but also achieves excellent re-identification performance.

I. Qualitative Results

To illustrate the advantages of our RDE, some retrieved examples for TIReID are presented in Figure 13. These results are obtained by testing the model trained on the CUHK-PEDES dataset with 20% NCs. From the examples, one can see that our RDE obtains more accurate and reasonable re-identification results. Simultaneously, in some inaccurate results (*e.g.*, the results (b) and (d)) obtained by IRRA, we find that the visual information of the retrieved image often only matches part of the text query, which indicates that the model cannot learn complete alignment knowledge. We think the reason is that the NCs mislead the model of IRRA to focus on some wrong visual-semantic associations. In contrast, our RDE could filter out erroneous correspondences to learn reliable and accurate cross-modal knowledge, thus achieving high robustness and better results.

References

- [1] Devansh Arpit, Stanisław Jastrzębski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S. Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, and Simon Lacoste-Julien. A closer look at memorization in deep networks. In *Proceedings of the 34th International Conference on Machine Learning*, pages 233–242. PMLR, 2017. 2, 3, 4
- [2] Yang Bai, Min Cao, Daming Gao, Ziqiang Cao, Chen Chen, Zhenfeng Fan, Liqiang Nie, and Min Zhang. Rasa: relation and sensitivity aware representation learning for text-based person search. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence*, pages 555–563, 2023. 2, 12
- [3] Yang Bai, Jingyao Wang, Min Cao, Chen Chen, Ziqiang Cao, Liqiang Nie, and Min Zhang. Text-based person search without parallel image-text data. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 757–767, 2023.
- [4] Min Cao, Yang Bai, Ziyin Zeng, Mang Ye, and Min Zhang. An empirical study of clip for text-based person search. *arXiv preprint arXiv:2308.10045*, 2023. 2, 12
- [5] Yuhao Chen, Guoqing Zhang, Yujiang Lu, Zhenxing Wang, and Yuhui Zheng. Tipcb: A simple but effective part-based convolutional baseline for text-based person search. *Neuro-computing*, 494:171–181, 2022. 12
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 2
- [7] Haiwen Diao, Ying Zhang, Lin Ma, and Huchuan Lu. Similarity reasoning and filtration for image-text matching. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1218–1226, 2021. 5
- [8] Zefeng Ding, Changxing Ding, Zhiyin Shao, and Dacheng Tao. Semantically self-aligned network for text-to-image part-aware person re-identification. *arXiv preprint arXiv:2107.12666*, 2021. 6, 7, 10, 12
- [9] Jianfeng Dong, Xirong Li, Chaoxi Xu, Xun Yang, Gang Yang, Xun Wang, and Meng Wang. Dual encoding for video retrieval by text. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(8):4065–4080, 2021. 5
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2
- [11] Chanhon Eom and Bumsub Ham. Learning disentangled representation for robust person re-identification. *Advances in neural information processing systems*, 32, 2019. 1
- [12] Fartash Faghri, David J Fleet, Jamie Ryan Kiros, and Sanja Fidler. Vse++: Improving visual-semantic embeddings with hard negatives. *arXiv preprint arXiv:1707.05612*, 2017. 2, 3
- [13] Ammarah Farooq, Muhammad Awais, Josef Kittler, and Syed Safwan Khalid. Axm-net: Implicit cross-modal feature alignment for person re-identification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 4477–4485, 2022. 12
- [14] Yanglin Feng, Hongyuan Zhu, Dezhong Peng, Xi Peng, and Peng Hu. Rono: Robust discriminative learning with noisy labels for 2d-3d cross-modal retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11610–11619, 2023. 2
- [15] Takuro Fujii and Shuhei Tarashima. Bilma: Bidirectional local-matching for text-based person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2786–2790, 2023. 12
- [16] Chenyang Gao, Guanyu Cai, Xinyang Jiang, Feng Zheng, Jun Zhang, Yifei Gong, Pai Peng, Xiaowei Guo, and Xing Sun. Contextual non-local alignment over full-scale rep-

- resentation for text-based person search. *arXiv preprint arXiv:2101.03036*, 2021. 2, 12
- [17] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. *Advances in neural information processing systems*, 31, 2018. 4
- [18] Haochen Han, Kaiyao Miao, Qinghua Zheng, and Minnan Luo. Noisy correspondence learning with meta similarity correction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7517–7526, 2023. 3
- [19] Xiao Han, Sen He, Li Zhang, and Tao Xiang. Text-based person search with limited data. *arXiv preprint arXiv:2110.10807*, 2021. 1, 2, 12
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4
- [21] Peng Hu, Xi Peng, Hongyuan Zhu, Liangli Zhen, and Jie Lin. Learning cross-modal retrieval with noisy labels. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5403–5413, 2021. 2
- [22] Peng Hu, Zhenyu Huang, Dezhong Peng, Xu Wang, and Xi Peng. Cross-modal retrieval with partially mismatched pairs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–15, 2023. 3
- [23] Zhenyu Huang, Guocheng Niu, Xiao Liu, Wenbiao Ding, Xinyan Xiao, Hua Wu, and Xi Peng. Learning with noisy correspondence for cross-modal matching. *Advances in Neural Information Processing Systems*, 34:29406–29419, 2021. 2, 3, 4, 5
- [24] Ding Jiang and Mang Ye. Cross-modal implicit relation reasoning and aligning for text-to-image person retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2787–2797, 2023. 1, 2, 3, 4, 6, 7, 8, 10, 11, 12
- [25] Ya Jing, Chenyang Si, Junbo Wang, Wei Wang, Liang Wang, and Tieniu Tan. Pose-guided multi-granularity attention network for text-based person search. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 11189–11196, 2020. 1
- [26] Junnan Li, Richard Socher, and Steven CH Hoi. Dividemix: Learning with noisy labels as semi-supervised learning. *arXiv preprint arXiv:2002.07394*, 2020. 2, 4, 5
- [27] Shuang Li, Tong Xiao, Hongsheng Li, Bolei Zhou, Dayu Yue, and Xiaogang Wang. Person search with natural language description. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1970–1979, 2017. 1, 2, 6, 10
- [28] Shiping Li, Min Cao, and Min Zhang. Learning semantic-aligned feature representation for text-based person search. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2724–2728. IEEE, 2022. 12
- [29] Shenshen Li, Xing Xu, Yang Yang, Fumin Shen, Yijun Mo, Yujie Li, and Heng Tao Shen. Dcel: Deep cross-modal evidential learning for text-based person retrieval. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 6292–6300, 2023. 2, 12
- [30] Zheng Li, Caili Guo, Xin Wang, Zerun Feng, and Zhongtian Du. Selectively hard negative mining for alleviating gradient vanishing in image-text matching. *arXiv preprint arXiv:2303.00181*, 2023. 10
- [31] Yijie Lin, Jie Zhang, Zhenyu Huang, Jia Liu, Zujie Wen, and Xi Peng. Multi-granularity correspondence learning from long-term noisy videos. *arXiv preprint arXiv:2401.16702*, 2024. 2
- [32] Yangdi Lu, Yang Bo, and Wenbo He. An ensemble model for combating label noise. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*, pages 608–617, 2022. 2
- [33] Yiwei Ma, Xiaoshuai Sun, Jiayi Ji, Guannan Jiang, Weilin Zhuang, and Rongrong Ji. Beat: Bi-directional one-to-many embedding alignment for text-based person retrieval. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 4157–4168, 2023. 2, 12
- [34] Kai Niu, Yan Huang, Wanli Ouyang, and Liang Wang. Improving description-based person re-identification by multi-granularity image-text alignments. *IEEE Transactions on Image Processing*, 29:5542–5556, 2020. 1
- [35] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 11
- [36] Yang Qin, Dezhong Peng, Xi Peng, Xu Wang, and Peng Hu. Deep evidential learning with noisy correspondence for cross-modal retrieval. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 4948–4956, 2022. 2, 3, 7, 11
- [37] Yalan Qin, Nan Pu, and Hanzhou Wu. Edmc: Efficient multi-view clustering via cluster and instance space learning. *IEEE Transactions on Multimedia*, 2023. 2
- [38] Yalan Qin, Nan Pu, and Hanzhou Wu. Elastic multi-view subspace clustering with pairwise and high-order correlations. *IEEE Transactions on Knowledge and Data Engineering*, 2023. 2
- [39] Yang Qin, Yuan Sun, Dezhong Peng, Joey Tianyi Zhou, Xi Peng, and Peng Hu. Cross-modal active complementary learning with self-refining correspondence. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 3
- [40] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2, 7
- [41] Nikolaos Sarafianos, Xiang Xu, and Ioannis A Kakadiaris. Adversarial representation learning for text-to-image matching. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5814–5824, 2019. 12
- [42] Zhiyin Shao, Xinyu Zhang, Meng Fang, Zhifeng Lin, Jian Wang, and Changxing Ding. Learning granularity-unified representations for text-to-image person re-identification. In

- Proceedings of the 30th ACM International Conference on Multimedia*, 2022. 1, 2, 12
- [43] Zhiyin Shao, Xinyu Zhang, Changxing Ding, Jian Wang, and Jingdong Wang. Unified pre-training with pseudo texts for text-to-image person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11174–11184, 2023. 2, 12
 - [44] Fei Shen, Xiangbo Shu, Xiaoyu Du, and Jinhui Tang. Pedestrian-specific bipartite-aware similarity learning for text-based person retrieval. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 8922–8931, 2023. 2, 12
 - [45] Xiujun Shu, Wei Wen, Haoqian Wu, Keyu Chen, Yiran Song, Ruizhi Qiao, Bo Ren, and Xiao Wang. See finer, see more: Implicit modality alignment for text-based person retrieval. In *European Conference on Computer Vision*, pages 624–641. Springer, 2022. 1, 2, 7, 10, 11, 12
 - [46] Chengji Wang, Zhiming Luo, Yaojin Lin, and Shaozi Li. Text-based person search via multi-granularity embedding learning. In *IJCAI*, pages 1068–1074, 2021. 1, 2
 - [47] Zheng Wang, Ruimin Hu, Yi Yu, Chao Liang, and Wenxin Huang. Multi-level fusion for person re-identification with incomplete marks. In *Proceedings of the 23rd ACM international conference on Multimedia*, pages 1267–1270, 2015. 1
 - [48] Zhe Wang, Zhiyuan Fang, Jun Wang, and Yezhou Yang. Vi-taa: Visual-textual attributes alignment in person search by natural language. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XII 16*, pages 402–420. Springer, 2020. 12
 - [49] Zijie Wang, Aichun Zhu, Jingyi Xue, Xili Wan, Chao Liu, Tian Wang, and Yifeng Li. Caibc: Capturing all-round information beyond color for text-based person retrieval. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 5314–5322, 2022. 12
 - [50] Zijie Wang, Aichun Zhu, Jingyi Xue, Xili Wan, Chao Liu, Tian Wang, and Yifeng Li. Look before you leap: Improving text-based person retrieval by learning a consistent cross-modal common manifold. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 1984–1992, 2022. 12
 - [51] Longhui Wei, Shiliang Zhang, Wen Gao, and Qi Tian. Person transfer gan to bridge domain gap for person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 79–88, 2018. 10
 - [52] Yushuang Wu, Zizheng Yan, Xiaoguang Han, Guanbin Li, Changqing Zou, and Shuguang Cui. Lapscore: language-guided person search via color reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1624–1633, 2021. 1, 12
 - [53] Shuanglin Yan, Neng Dong, Liyan Zhang, and Jinhui Tang. Clip-driven fine-grained text-image person re-identification. *arXiv preprint arXiv:2210.10276*, 2022. 1, 2, 4, 7, 12
 - [54] Shuanglin Yan, Hao Tang, Liyan Zhang, and Jinhui Tang. Image-specific information suppression and implicit local alignment for text-based person search. *arXiv preprint arXiv:2208.14365*, 2022. 12
 - [55] Shuanglin Yan, Neng Dong, Jun Liu, Liyan Zhang, and Jinhui Tang. Learning comprehensive representations with richer self for text-to-image person re-identification. In *Proceedings of the 31st ACM international conference on multimedia*, pages 6202–6211, 2023. 2, 12
 - [56] Mouxing Yang, Zhenyu Huang, Peng Hu, Taihao Li, Jiancheng Lv, and Xi Peng. Learning with twin noisy labels for visible-infrared person re-identification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14308–14317, 2022. 2, 3
 - [57] Mouxing Yang, Yunfan Li, Peng Hu, Jinfeng Bai, Jiancheng Lv, and Xi Peng. Robust multi-view clustering with incomplete information. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(1):1055–1069, 2022. 2
 - [58] Mouxing Yang, Zhenyu Huang, and Xi Peng. Robust object re-identification with coupled noisy labels. *International Journal of Computer Vision*, pages 1–19, 2024. 2
 - [59] Huaiwen Zhang, Yang Yang, Fan Qi, Shengsheng Qian, and Changsheng Xu. Robust video-text retrieval via noisy pair calibration. *IEEE Transactions on Multimedia*, 2023. 2, 3
 - [60] Ying Zhang and Huchuan Lu. Deep cross-modal projection learning for image-text matching. In *Proceedings of the European conference on computer vision (ECCV)*, pages 686–701, 2018. 1, 2, 12
 - [61] Zhedong Zheng, Liang Zheng, Michael Garrett, Yi Yang, Mingliang Xu, and Yi-Dong Shen. Dual-path convolutional image-text embeddings with instance loss. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 16(2):1–23, 2020. 2, 12
 - [62] Aichun Zhu, Zijie Wang, Yifeng Li, Xili Wan, Jing Jin, Tian Wang, Fangqiang Hu, and Gang Hua. Dssl: Deep surroundings-person separation learning for text-based person retrieval. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 209–217, 2021. 1, 2, 6, 10, 12
 - [63] Haowei Zhu, Wenjing Ke, Dong Li, Ji Liu, Lu Tian, and Yi Shan. Dual cross-attention learning for fine-grained visual categorization and object re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4692–4702, 2022. 4