

# Hierarchical Identity Learning for Unsupervised Visible-Infrared Person Re-Identification

Haonan Shi, Yubin Wang, De Cheng, Lingfeng He, Nannan Wang, *Senior Member, IEEE*,  
Xinbo Gao, *Fellow, IEEE*

**Abstract**—Unsupervised visible-infrared person re-identification (USVI-ReID) aims to learn modality-invariant image features from unlabeled cross-modal person datasets by reducing the modality gap while minimizing reliance on costly manual annotations. Existing methods typically address USVI-ReID using cluster-based contrastive learning, which represents a person by a single cluster center. However, they primarily focus on the commonality of images within each cluster while neglecting the finer-grained differences among them. To address the limitation, we propose a Hierarchical Identity Learning (HIL) framework. Since each cluster may contain several smaller sub-clusters that reflect fine-grained variations among images, we generate multiple memories for each existing coarse-grained cluster via a secondary clustering. Additionally, we propose Multi-Center Contrastive Learning (MCCL) to refine representations for enhancing intra-modal clustering and minimizing cross-modal discrepancies. To further improve cross-modal matching quality, we design a Bidirectional Reverse Selection Transmission (BRST) mechanism, which establishes reliable cross-modal correspondences by performing bidirectional matching of pseudo-labels. Extensive experiments conducted on the SYSU-MM01 and RegDB datasets demonstrate that the proposed method outperforms existing approaches. The source code is available at: <https://github.com/haonanshi0125/HIL>.

**Index Terms**—Unsupervised visible-infrared person re-identification, Hierarchical learning, Cross-Modal matching, Clustering algorithm

## I. INTRODUCTION

VISIBLE-infrared person re-identification (VI-ReID) [1], [2], [3], [4], [5], [6] is an important research direction in the field of computer vision, aiming to match the images of the same person between the visible and infrared modalities. Compared to the extensively researched single-modality person ReID [7], [8], VI-ReID is a more challenging task due to the large modality gap between visible and infrared images. However, annotating cross-modal datasets demands more resources than single-modal datasets. To tackle the challenge of heavy annotations on large-scale cross-modal data, several semi-supervised methods [4], [9], [10] have been proposed for visible-infrared person re-identification,

This work was supported in part by the National Natural Science Foundation of China under Grants 62176198, U22A2096, in part by the Key R&D Program of Shaanxi Province under Grant 2024GX-YBXM135.

Haonan Shi, De Cheng, Lingfeng He, Nannan Wang and Xinbo Gao are with the State Key Laboratory of Integrated Services Networks, School of Telecommunications Engineering, Xidian University, Xi'an 710071, Shanxi, P. R. China (email: dcheng@xidian.edu.cn, lfhe@stu.xidian.edu.cn, nnwang@xidian.edu.cn and xbgao@mail.xidian.edu.cn).

Yubin Wang is with the Department of Computer Science and Technology, Tongji University, Shanghai 201804, China (e-mail: wangyubin2018@tongji.edu.cn).

(Corresponding author: De Cheng).

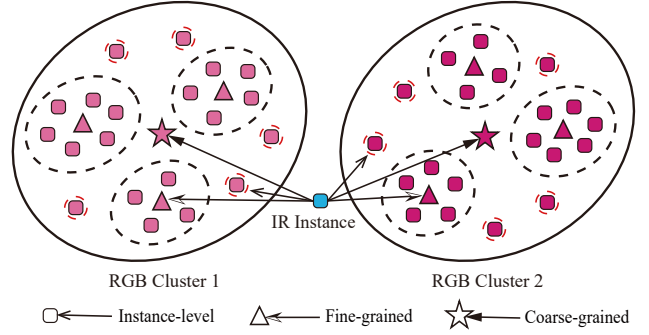


Fig. 1. Illustration of the hierarchical structure of feature alignment. Hierarchical identity information is extracted at three levels through two rounds of clustering, facilitating the refinement of representations within and across modalities via aligning instances with other instances, fine-grained centers, and coarse-grained centers.

leveraging both labeled and unlabeled data to learn modality-invariant and identity-discriminative representations. Although these methods have shown promising results, they still depend on a certain amount of annotated data.

To address scenarios lacking annotated data, some unsupervised methods [11], [12], [13], [14], [15] have been proposed, which employ a single memory bank for each cluster in both visible and infrared modalities. During training, the memory centers act as prototypes, and contrastive loss between query images and these prototypes, which are aggregated based on their similarity, is minimized. While single-memory representation learning offers simplicity, it may fail to capture the diverse perspectives and intricate details of the same identity, potentially resulting in information loss. To overcome this limitation, MMM [16] employs multi-memory representations for identity generation, which enhances the capacity to capture fine-grained variations. However, simply relying on multi-memory representations without preserving coarse-grained information may lead to overfitting, as it neglects the macro-level commonality and informative features of the identity. As a result, the model tends to concentrate on a few specific fine-grained patterns, which can adversely affect optimization for modality-invariant features. As shown in Figure 1, we argue that incorporating a hierarchical structure of features, including instance-level, coarse-grained, and fine-grained features, enhances the refinement of representations and promotes better alignment both within and across modalities.

To this end, we propose a Hierarchical Identity Learning (HIL) framework, which collaboratively integrates both coarse-grained and fine-grained pseudo-labels. Specifically,

our framework leverages a two-stage clustering strategy to refine the representation of identity features. In the first stage, we employ DBSCAN [17] to generate coarse-grained pseudo-labels, grouping visible and infrared instance features into identity clusters. However, as each cluster may contain multiple smaller sub-clusters that reflect fine-grained variations among images, a single clustering stage may fail to capture the intricate details and diverse perspectives of the same identity. To overcome this limitation, we introduce a secondary clustering step within each coarse-grained cluster, which facilitates a more comprehensive understanding of identity features.

To effectively utilize and model the hierarchical identity information, we propose a Multi-Center Contrastive Learning (MCCL) strategy. By leveraging the centers of multi-granularity clusters as reference points, MCCL constructs robust positive and negative sample sets. This approach allows the model to refine representations through contrastive learning both within and across modalities, thereby enhancing intra-modal clustering and minimizing cross-modal discrepancies. Furthermore, we design a Bidirectional Reverse Selection Transmission (BRST) mechanism to improve cross-modal matching quality. This mechanism performs bidirectional matching of pseudo-labels between visible and infrared modalities, enabling the establishment of reliable cross-modal correspondences. A reverse selection process is integrated to filter out unreliable pseudo-label matches, thereby enhancing their overall quality and robustness.

By combining these components, our proposed HIL framework addresses key limitations and achieves improved performance in cross-modal retrieval tasks. Our contributions can be summarized as follows:

- We propose a Hierarchical Identity Learning (HIL) framework that employs a secondary clustering step within coarse-grained clusters to capture fine-grained identity variations, enabling robust representation.
- We propose a Multi-Center Contrastive Learning (MCCL) strategy that utilizes hierarchical identity information to construct robust positive and negative sample sets and refine representations for enhancing intra-modal clustering and minimizing cross-modal discrepancies.
- We design a Bidirectional Reverse Selection Transmission (BRST) mechanism that performs bidirectional matching of pseudo-labels between modalities, ensuring robust cross-modal correspondences by filtering unreliable matches through a reverse selection process.
- Extensive experiments conducted on the SYSU-MM01 and RegDB datasets demonstrate that the proposed method outperforms existing approaches in various settings.

## II. RELATED WORK

### A. Supervised Visible-Infrared Person ReID

Supervised visible-infrared person ReID methods primarily focus on bridging the gap between the two different modalities using labeled person images. MPANet [18] proposed a joint modality and pattern alignment network to discover cross-modality nuances. CAJ [19] proposed a channel augmented

joint learning strategy to improve the robustness against color variations by randomly exchanging the color channels. CIFT [20] proposed a counterfactual intervention feature transfer method to address the balance gap between training-test modality and suboptimal topology structure problems. FMCNet [21] employed GANs to compensate for missing modality-specific information at the feature level. TransVI [22] designed a transformer-based visible-infrared network with a two-stream structure to capture modality-specific features and learn shared knowledge. SEFL [23] proposed a shape-erased feature learning paradigm to eliminate body-shape-related information from the learned features. SAAI [24] proposed a semantic alignment and affinity inference framework to explore the joint application of semantic-aligned feature learning and the affinity inference method. IDKL [25] proposed an implicit discriminative knowledge learning network to uncover and leverage the implicit discriminative information contained within the modality-specific. These methods have effectively reduced the modality gap, demonstrating their usefulness for supervised VI-ReID. However, the performance of these methods requires extensive human-labeled cross-modal data, which is time-consuming and expensive.

### B. Unsupervised Single-Modality Person ReID

Unsupervised single-modality person ReID tasks endeavor to learn robust representations for unlabeled person images within a single modality. In order to mitigate the effects of noisy pseudo-labels, MMT [26] introduced a mutual mean-teaching framework that provides reliable soft pseudo-labels. SPCL [27] gradually generates more robust clusters through a self-paced contrastive learning framework with hybrid memory. ICE [28] introduced an inter-instance contrastive encoding method that enhances cluster compactness and improves pseudo-labels quality. RLCC [29] proposed a method to accurately estimate pseudo-label similarities between consecutive training generations using clustering consensus, and to refine pseudo-labels through temporal propagation and ensembling. IICS [30] proposed a two-stage similarity computation strategy that separately models intra-camera and inter-camera relations to generate more reliable pseudo-labels. To further improve the quality of the pseudo-labels, ISE [31] designed an implicit sample extension method to strengthen the reliability of clusters. PPLR [8] proposed a part-based pseudo-label refinement framework that reduces pseudo-label noise by utilizing the complementary relationship between global and part features. Cluster-Contrast [32] proposed a cluster contrast method that stores unique centroid features and performs contrastive learning at the cluster level. However, they are difficult to directly apply to solving unsupervised visible-infrared person ReID tasks due to the large modality gap.

### C. Unsupervised Visible-Infrared Person ReID

Unsupervised visible-infrared person ReID tasks aim to learn modality-invariant representations and establish reliable cross-modal associations between visible and infrared modalities without identity annotations. ADCA [5] proposed an augmented dual-contrastive aggregation learning framework

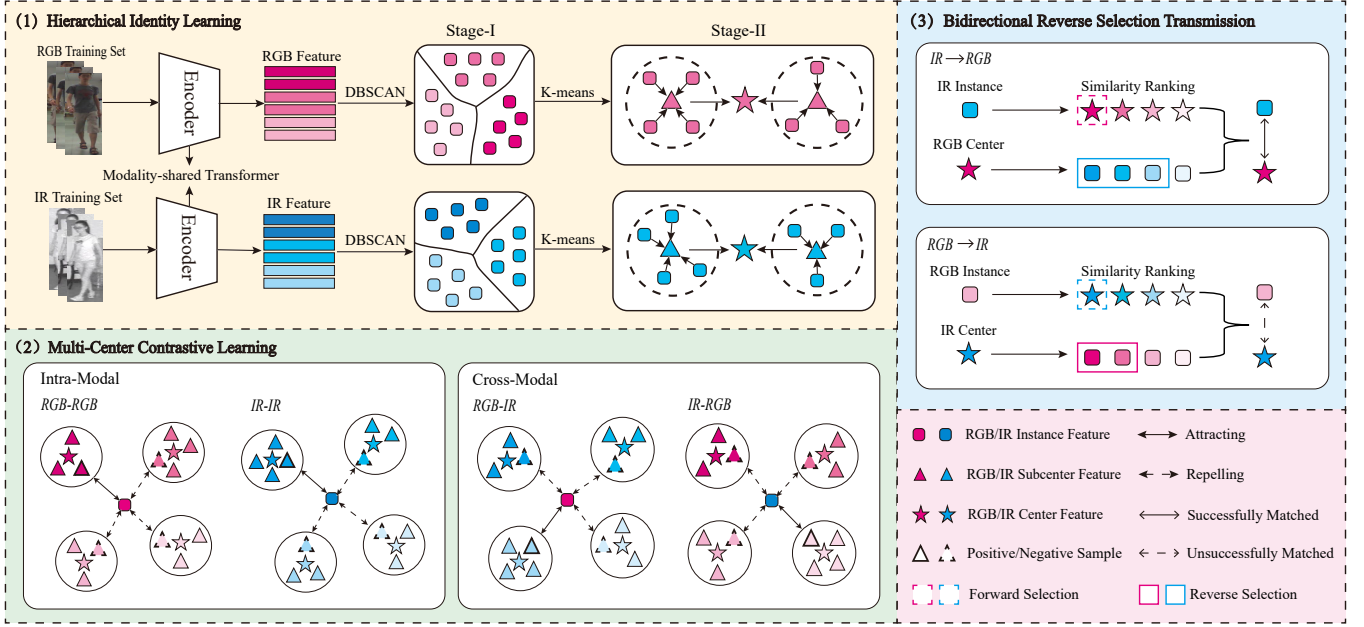


Fig. 2. Illustration of the Hierarchical Identity Learning (HIL) framework. Given unlabeled visible-infrared data, this framework generates more comprehensive hierarchical identity information, with the Multi-Center Contrastive Learning (MCCL) strategy utilizing these centers to construct robust positive and negative sample sets for refining representations via contrastive learning. For the cross-modal alignment, the Bidirectional Reverse Selection Transmission (BRST) mechanism performs bidirectional matching of pseudo-labels between two modalities.

based on the ideas of homogeneous joint learning and heterogeneous aggregation. CHCR [12] introduced a cross-modality hierarchical clustering and refinement method by enhancing modality-invariant feature learning and strengthening the reliability of pseudo-labels. DOTLA [33] employed the optimal transport strategy to assign pseudo-labels from one modality to another modality at the instance level. PGM [34] proposed a progressive graph matching method to establish reliable cross-modal correspondences. SDCL [35] developed a shallow-deep collaborative learning framework based on the transformer architecture, which reduces the modality gap through the collaboration of shallow and deep features. PCLHD [36] introduced a method of progressive contrastive learning with hard and dynamic prototypes, effectively learning commonality, divergence, and variety. MMM [16] proposed a multi-memory matching framework to effectively capture intra-class nuances and establish reliable cross-modality correspondences. RPNR [37] introduced a robust pseudo-label learning with neighbor relation framework, enhancing pseudo-label reliability and strengthening cross-modality alignment. Although these approaches demonstrate satisfactory performance, they are insufficient for effectively learning modality-invariant features without the simultaneous integration of multi-level identity information.

### III. METHOD

#### A. Problem Formulation and Overview

Given an unlabeled visible-infrared ReID dataset  $\mathcal{X} = \{\mathcal{X}^v, \mathcal{X}^r\}$ , where  $\mathcal{X}^v = \{\mathbf{x}_k^v\}_{k=1}^{N_v}$  denotes the visible dataset containing  $N_v$  unlabeled samples and  $\mathcal{X}^r = \{\mathbf{x}_k^r\}_{k=1}^{N_r}$  denotes the infrared dataset containing  $N_r$  unlabeled samples. Our

objective is to train a deep neural network  $f_\theta(\cdot)$  that maps images from both modalities into a shared feature space. This shared feature space is designed to produce  $d$ -dimensional representations such that  $\mathbf{f}_k^v = f_\theta(\mathbf{x}_k^v) \in \mathbb{R}^d$  and  $\mathbf{f}_k^r = f_\theta(\mathbf{x}_k^r) \in \mathbb{R}^d$ . Following the mapping process, we can obtain the instance feature set  $\mathcal{F} = \{\mathcal{F}^v, \mathcal{F}^r\}$ , where  $\mathcal{F}^v = \{\mathbf{f}_k^v\}_{k=1}^{N_v}$  denotes the  $N_v$  visible instance features and  $\mathcal{F}^r = \{\mathbf{f}_k^r\}_{k=1}^{N_r}$  denotes the  $N_r$  infrared instance features. By ensuring the representations are modality-invariant, the network enables effective matching of images of the same identity across the two modalities.

Our proposed Hierarchical Identity Learning (HIL) framework, illustrated in Figure 2, integrates the Multi-Center Contrastive Learning (MCCL) strategy and the Bidirectional Reverse Selection Transmission (BRST) mechanism to enhance identity representation learning. This framework is built upon a dual-path transformer architecture inspired by [35]. Within this architecture, instance features are clustered in two stages to produce coarse-grained and fine-grained pseudo-labels. Based on the identity information, MCCL refines representations via contrastive learning on multiple selected centers, aiming to derive robust and reliable intra-modal representations and minimize the modality gap. Additionally, BRST refines cross-modal matching by establishing consistent cross-modal correspondences through bidirectional pseudo-label matching in a reverse manner. By integrating hierarchical identity learning with the cross-modal feature matching algorithm, HIL effectively addresses the challenge of feature robustness and identity alignment across modalities.

## B. Hierarchical Identity Learning

a) *Coarse-grained Pseudo-Label Generation*: To obtain pseudo-labels in the unsupervised setting, we employ DBSCAN [17] for clustering visible and infrared instance features into  $M_v$  and  $M_r$  groups, respectively, at the first stage. After clustering, we can obtain the coarse-grained pseudo-label set  $\mathcal{Y}^c = \{\mathcal{Y}^{vc}, \mathcal{Y}^{rc}\}$ , where  $\mathcal{Y}^{vc} = \{\mathbf{y}_i^{vc}\}_{i=1}^{M_v}$  represents the  $M_v$  visible coarse-grained pseudo-labels and  $\mathcal{Y}^{rc} = \{\mathbf{y}_i^{rc}\}_{i=1}^{M_r}$  represents the  $M_r$  infrared coarse-grained pseudo-labels. By calculating the mean feature of all instance features within each corresponding coarse-grained cluster, we can obtain a set of coarse-grained cluster features  $\mathcal{U}^c = \{\mathcal{U}^{vc}, \mathcal{U}^{rc}\}$ , where  $\mathcal{U}^{vc} = \{\mathbf{u}_i^{vc}\}_{i=1}^{M_v}$  represents the  $M_v$  visible coarse-grained cluster features and  $\mathcal{U}^{rc} = \{\mathbf{u}_i^{rc}\}_{i=1}^{M_r}$  represents the  $M_r$  infrared coarse-grained cluster features. The coarse-grained cluster features are calculated as follows:

$$\mathbf{u}_i^{vc} = \frac{1}{|\mathcal{H}_i^{vc}|} \sum_{\mathbf{f}_n^{vc} \in \mathcal{H}_i^{vc}} \mathbf{f}_n^{vc}, \quad (1)$$

$$\mathbf{u}_i^{rc} = \frac{1}{|\mathcal{H}_i^{rc}|} \sum_{\mathbf{f}_n^{rc} \in \mathcal{H}_i^{rc}} \mathbf{f}_n^{rc}, \quad (2)$$

where  $\mathbf{f}_n^{vc}$  and  $\mathbf{f}_n^{rc}$  are visible and infrared instance features within the same coarse-grained pseudo-label, respectively.  $\mathcal{H}_i^{vc}$  and  $\mathcal{H}_i^{rc}$  are visible and infrared coarse-grained cluster sets, respectively. The operator  $|\cdot|$  counts the number of samples of a set.

During each training iteration, the coarse-grained cluster features of both modalities are updated using a momentum-based strategy:

$$\mathbf{u}_{i,t}^{vc} \leftarrow \alpha \mathbf{u}_{i,t-1}^{vc} + (1 - \alpha) q^v, \quad q^v \in \mathcal{H}_i^{vc}, \quad (3)$$

$$\mathbf{u}_{i,t}^{rc} \leftarrow \alpha \mathbf{u}_{i,t-1}^{rc} + (1 - \alpha) q^r, \quad q^r \in \mathcal{H}_i^{rc}, \quad (4)$$

where  $q^v$  and  $q^r$  are visible and infrared query features. The  $\alpha$  is a momentum hyperparameter.  $t$  and  $t - 1$  refer to the current and last iterations, respectively.

Given visible and infrared query features  $q^v$  and  $q^r$ , we compute the contrastive loss for visible and infrared modalities by the following equations:

$$\mathcal{L}_{id}^v = -\log \frac{\exp(q^v \cdot \mathbf{u}_+^{vc}/\tau)}{\sum_{i=0}^{M_v} \exp(q^v \cdot \mathbf{u}_i^{vc}/\tau)}, \quad (5)$$

$$\mathcal{L}_{id}^r = -\log \frac{\exp(q^r \cdot \mathbf{u}_+^{rc}/\tau)}{\sum_{i=0}^{M_r} \exp(q^r \cdot \mathbf{u}_i^{rc}/\tau)}, \quad (6)$$

where  $\mathbf{u}_+^{vc}$  and  $\mathbf{u}_+^{rc}$  are the positive coarse-grained cluster features corresponding to the pseudo-labels of  $q^v$  and  $q^r$ , respectively. The  $\tau$  is a temperature hyperparameter. The overall identity loss is defined as:

$$\mathcal{L}_{id} = \mathcal{L}_{id}^v + \mathcal{L}_{id}^r. \quad (7)$$

b) *Fine-grained Pseudo-Label Generation*: After the first stage, the clustering process may not capture the diverse perspectives and intricate details associated with the same identity. To further explore fine-grained variations within existing identity clusters, we apply a secondary clustering algorithm, such as K-means [38]. This secondary clustering step is performed within each coarse-grained cluster to refine the identity representation. As a result, we can obtain a fine-grained pseudo-label set  $\mathcal{Y}^f = \{\mathcal{Y}^{vf}, \mathcal{Y}^{rf}\}$ , where  $\mathcal{Y}^{vf} = \{\mathbf{y}_{ij}^{vf}\}_{i=1, j=1}^{i=M_v, j=K}$  represents the visible fine-grained pseudo-label set and  $\mathcal{Y}^{rf} = \{\mathbf{y}_{ij}^{rf}\}_{i=1, j=1}^{i=M_r, j=K}$  represents the infrared fine-grained pseudo-label set. Specifically,  $K$  is the number of fine-grained clusters within each coarse-grained cluster. By calculating the mean feature of all instance features within each corresponding fine-grained cluster, we can obtain a fine-grained cluster feature set  $\mathcal{U}^f = \{\mathcal{U}^{vf}, \mathcal{U}^{rf}\}$ , where  $\mathcal{U}^{vf} = \{\mathbf{u}_{ij}^{vf}\}_{i=1, j=1}^{i=M_v, j=K}$  represents the visible fine-grained cluster features and  $\mathcal{U}^{rf} = \{\mathbf{u}_{ij}^{rf}\}_{i=1, j=1}^{i=M_r, j=K}$  represents the infrared fine-grained cluster features. The fine-grained cluster features are calculated as follows:

$$\mathbf{u}_{ij}^{vf} = \frac{1}{|\mathcal{H}_{ij}^{vf}|} \sum_{\mathbf{f}_n^{vf} \in \mathcal{H}_{ij}^{vf}} \mathbf{f}_n^{vf}, \quad (8)$$

$$\mathbf{u}_{ij}^{rf} = \frac{1}{|\mathcal{H}_{ij}^{rf}|} \sum_{\mathbf{f}_n^{rf} \in \mathcal{H}_{ij}^{rf}} \mathbf{f}_n^{rf}, \quad (9)$$

where  $\mathbf{f}_n^{vf}$  and  $\mathbf{f}_n^{rf}$  are visible and infrared instance features within the same fine-grained pseudo-label, respectively.  $\mathcal{H}_{ij}^{vf}$  and  $\mathcal{H}_{ij}^{rf}$  are visible and infrared fine-grained cluster sets, respectively. The operator  $|\cdot|$  counts the number of samples in a set.

The hierarchical identity learning framework enables a comprehensive analysis of identities and facilitates the convergence of instance features toward their corresponding explicit centers. This process enhances the formation of positive and negative pairs, which are critical for the subsequent multi-center contrastive learning strategy.

## C. Neighbor Contrastive Learning

Based on the obtained hierarchical identity information, we further expand the contrastive learning by exploring the relationship between different instances. For each visible instance feature, we calculate its cosine similarity with all visible instance features to find the one with the highest similarity. The cosine similarity between the instance feature  $\mathbf{f}_k^v$  and the instance feature  $\mathbf{f}_n^v$  is computed as follows:

$$\text{Sim}(\mathbf{f}_k^v, \mathbf{f}_n^v) = \frac{\mathbf{f}_k^v \cdot \mathbf{f}_n^v}{\|\mathbf{f}_k^v\| \|\mathbf{f}_n^v\|}. \quad (10)$$

The instance feature with the highest similarity for a given instance feature  $\mathbf{f}_k^v$  can be identified as:

$$\tilde{\mathbf{f}}_k^v = \arg \max_{\mathbf{f}_n^v \in \mathcal{F}^v} \text{Sim}(\mathbf{f}_k^v, \mathbf{f}_n^v). \quad (11)$$

Next, we select reliable neighbors as the positive sample set  $\mathcal{R}_k^v$  and choose the remaining visible instance features as the negative sample set  $\mathcal{W}_k^v$  for the instance feature  $\mathbf{f}_k^v$ . The  $\mathcal{R}_k^v$  and  $\mathcal{W}_k^v$  can be obtained by:

$$\mathcal{R}_k^v = \left\{ \mathbf{f}_n^v \mid \text{Sim}(\mathbf{f}_k^v, \mathbf{f}_n^v) > \beta \cdot \text{Sim}(\mathbf{f}_k^v, \tilde{\mathbf{f}}_k^v) \right\}, \quad (12)$$

$$\mathcal{W}_k^v = \left\{ \mathbf{f}_n^v \mid \text{Sim}(\mathbf{f}_k^v, \mathbf{f}_n^v) \leq \beta \cdot \text{Sim}(\mathbf{f}_k^v, \tilde{\mathbf{f}}_k^v) \right\}, \quad (13)$$

where  $\beta$  is a selection threshold hyperparameter.

By obtaining the positive and negative sample sets, the contrastive loss for visible-visible  $\mathcal{L}_{neighbor}^{vv}$  is defined as follows:

$$S_{k,r}^v = \sum_{\mathbf{f} \in \mathcal{R}_k^v} \exp \left( \frac{\text{Sim}(\mathbf{f}_k^v, \mathbf{f})}{\tau} \right), \quad (14)$$

$$S_{k,w}^v = \sum_{\mathbf{f} \in \mathcal{W}_k^v} \exp \left( \frac{\text{Sim}(\mathbf{f}_k^v, \mathbf{f})}{\tau} \right), \quad (15)$$

$$\mathcal{L}_{neighbor}^{vv} = -\frac{1}{N_v} \sum_{k=1}^{N_v} \log \frac{S_{k,r}^v}{S_{k,r}^v + S_{k,w}^v}, \quad (16)$$

where  $\tau$  is a temperature hyperparameter that scales the similarities. By minimizing this contrastive loss, the model learns to maximize the similarities between the instance features and their positive samples while minimizing the similarities with their negative samples.

Similarly, the contrastive loss for infrared-infrared  $\mathcal{L}_{neighbor}^{rr}$ , infrared-visible  $\mathcal{L}_{neighbor}^{rv}$ , and visible-infrared  $\mathcal{L}_{neighbor}^{vr}$  can be obtained by similar ways. The final optimization for neighbor contrastive learning is denoted by the following combination:

$$\mathcal{L}_{neighbor} = \mathcal{L}_{neighbor}^{vv} + \mathcal{L}_{neighbor}^{rr} + \mathcal{L}_{neighbor}^{rv} + \mathcal{L}_{neighbor}^{vr}. \quad (17)$$

### D. Multi-Center Contrastive Learning

Since hierarchical identity information has been obtained, it is intuitive to utilize the centers of multi-granularity clusters as reference points for constructing positive and negative sample sets in contrastive learning. For each visible instance feature, we calculate its cosine similarity with all visible fine-grained cluster features to find the one with the highest similarity. The cosine similarity between the instance feature  $\mathbf{f}_k^v$  and the fine-grained cluster feature  $\mathbf{u}_{ij}^{vf}$  is computed as follows:

$$\text{Sim}(\mathbf{f}_k^v, \mathbf{u}_{ij}^{vf}) = \frac{\mathbf{f}_k^v \cdot \mathbf{u}_{ij}^{vf}}{\|\mathbf{f}_k^v\| \|\mathbf{u}_{ij}^{vf}\|}. \quad (18)$$

The fine-grained cluster feature with the highest similarity for a given instance feature  $\mathbf{f}_k^v$  can be identified as:

$$\tilde{\mathbf{u}}_k^v = \arg \max_{\mathbf{u}_{ij}^{vf} \in \mathcal{U}^{vf}} \text{Sim}(\mathbf{f}_k^v, \mathbf{u}_{ij}^{vf}). \quad (19)$$

Next, we select all fine-grained cluster features within the coarse-grained cluster of  $\tilde{\mathbf{u}}_k^v$  as the positive sample set  $\mathcal{P}_k^v = \{\mathbf{u}_{sj}^{vf}\}_{j=1}^K$  for the instance feature  $\mathbf{f}_k^v$ . Here,  $s$  denotes

the coarse-grained pseudo-label of the selected fine-grained feature  $\tilde{\mathbf{u}}_k^v$ . For the remaining visible coarse-grained clusters, we choose the most similar fine-grained cluster feature within each coarse-grained cluster as the negative sample set  $\mathcal{N}_k^v$  for the instance feature  $\mathbf{f}_k^v$ , which can be represented as:

$$\mathcal{N}_k^v = \bigcup_{1 \leq i \leq M_v \wedge i \neq s} \left\{ \arg \max_{\mathbf{u}_{ij}^{vf} \in \mathcal{U}^{vf}} \text{Sim}(\mathbf{f}_k^v, \mathbf{u}_{ij}^{vf}) \right\}. \quad (20)$$

By obtaining the positive and negative sample sets, the contrastive loss for visible-visible  $\mathcal{L}_{mccl}^{vv}$  is defined as follows:

$$S_{k,p}^v = \sum_{\mathbf{u} \in \mathcal{P}_k^v} \exp \left( \frac{\text{Sim}(\mathbf{f}_k^v, \mathbf{u})}{\tau} \right), \quad (21)$$

$$S_{k,n}^v = \sum_{\mathbf{u} \in \mathcal{N}_k^v} \exp \left( \frac{\text{Sim}(\mathbf{f}_k^v, \mathbf{u})}{\tau} \right), \quad (22)$$

$$\mathcal{L}_{mccl}^{vv} = -\frac{1}{N_v} \sum_{k=1}^{N_v} \log \frac{S_{k,p}^v}{S_{k,p}^v + S_{k,n}^v}, \quad (23)$$

where  $\tau$  is a temperature hyperparameter that scales the similarities. By minimizing this contrastive loss, the model learns to maximize the similarities between the instance features and their positive samples while minimizing the similarities with their negative samples.

Similarly, the contrastive loss for infrared-infrared  $\mathcal{L}_{mccl}^{rr}$ , infrared-visible  $\mathcal{L}_{mccl}^{rv}$ , and visible-infrared  $\mathcal{L}_{mccl}^{vr}$  can be obtained by similar ways. The final optimization for multi-center contrastive learning is denoted by the following combination:

$$\mathcal{L}_{mccl} = \mathcal{L}_{mccl}^{vv} + \mathcal{L}_{mccl}^{rr} + \mathcal{L}_{mccl}^{rv} + \mathcal{L}_{mccl}^{vr}. \quad (24)$$

The overall loss is denoted as:

$$\mathcal{L}_{total} = \mathcal{L}_{id} + \lambda_1 \mathcal{L}_{neighbor} + \lambda_2 \mathcal{L}_{mccl}, \quad (25)$$

where  $\lambda_1$  and  $\lambda_2$  are hyperparameters to balance the loss terms. Through this optimization, we enable the model to refine representations, thereby improving intra-modal clustering and minimizing cross-modal discrepancies. The training process in one epoch is illustrated in Algorithm 1.

### E. Bidirectional Reverse Selection Transmission

We introduce the Bidirectional Reverse Selection Transmission (BRST) mechanism to ensure robust and accurate label assignment for visible instance features by utilizing the coarse-grained cluster features of the infrared modality, and vice versa. These features serve as comprehensive representations in the hierarchical structure by encapsulating fine-grained variations, making them effective for achieving reliable identity alignment across modalities.

Specifically, for each visible instance feature, we calculate its cosine similarity with all infrared coarse-grained cluster features. This computation yields a similarity matrix  $\mathbf{S} \in \mathbb{R}^{N_v \times M_r}$ , where each row corresponds to one of the  $N_v$  visible instances, and each column corresponds to one of the  $M_r$  infrared coarse-grained clusters. Each element  $\mathbf{S}_{ki}$  in the matrix represents the similarity score between the  $k$ -th visible instance feature and the  $i$ -th infrared cluster feature.

**Algorithm 1:** Training process in one epoch**Require:**

Network  $f_\theta$ ; Current epoch number  $epoch$ ;  
Iterations per epoch  $T_t$ .

**Input:**

Visible training dataset  $\mathcal{X}^v = \{\mathbf{x}_k^v\}_{k=1}^{N_v}$ ;  
Infrared training dataset  $\mathcal{X}^r = \{\mathbf{x}_k^r\}_{k=1}^{N_r}$ .

**Output:**

Updated network  $f_\theta$ .

- 1:** Extract visible features  $\mathcal{F}^v = \{\mathbf{f}_k^v\}_{k=1}^{N_v}$  and infrared features  $\mathcal{F}^r = \{\mathbf{f}_k^r\}_{k=1}^{N_r}$ ;
- 2:** Cluster features by DBSCAN algorithm to obtain the coarse-grained pseudo-labels  $\mathcal{Y}^{vc} = \{\mathbf{y}_i^{vc}\}_{i=1}^{M_v}$  and  $\mathcal{Y}^{rc} = \{\mathbf{y}_i^{rc}\}_{i=1}^{M_r}$ ;
- 3:** Initialize the coarse-grained cluster features  $\mathcal{U}^{vc} = \{\mathbf{u}_i^{vc}\}_{i=1}^{M_v}$  and  $\mathcal{U}^{rc} = \{\mathbf{u}_i^{rc}\}_{i=1}^{M_r}$  by cluster centroids;
- 4:** Cluster features within the same coarse-grained pseudo-label by K-means algorithm to obtain the fine-grained pseudo-labels  $\mathcal{Y}^{vf} = \{\mathbf{y}_{ij}^{vf}\}_{i=1, j=1}^{i=M_v, j=K}$  and  $\mathcal{Y}^{rf} = \{\mathbf{y}_{ij}^{rf}\}_{i=1, j=1}^{i=M_r, j=K}$ ;
- 5:** Initialize the fine-grained cluster features  $\mathcal{U}^{vf} = \{\mathbf{u}_{ij}^{vf}\}_{i=1, j=1}^{i=M_v, j=K}$  and  $\mathcal{U}^{rf} = \{\mathbf{u}_{ij}^{rf}\}_{i=1, j=1}^{i=M_r, j=K}$  by cluster centroids;
- 6:** if  $epoch \% 2 == 0$  then
  - Obtain the modality-unified pseudo-labels  $\hat{\mathcal{Y}}^v = \{\hat{\mathbf{y}}_i^v\}_{i=1}^{M_v}$  by BRST mechanism;
- end**
- if**  $epoch \% 2 == 1$  **then**
  - Obtain the modality-unified pseudo-labels  $\hat{\mathcal{Y}}^r = \{\hat{\mathbf{y}}_i^r\}_{i=1}^{M_r}$  by BRST mechanism;
- end**
- 7:** for  $t$  in  $[1, T_t]$  **do**
  - Compute  $\mathcal{L}_{id}$  as Eq.7;
  - Compute  $\mathcal{L}_{neighbor}$  as Eq.17;
  - Compute  $\mathcal{L}_{mccl}$  as Eq.24;
  - Compute the total loss  $\mathcal{L}_{total}$  as Eq.25;
  - Update the coarse-grained cluster features;
  - Update the fine-grained cluster features;
  - Update the network parameters  $\theta$ .
- end**

To assign labels to visible instances, we first calculate the maximum similarity value and its corresponding column for each visible instance feature  $\mathbf{f}_k^v$ , which is defined as:

$$m_r = \max_i \mathbf{S}_{ki}, i^* = \arg \max_i \mathbf{S}_{ki}, \quad (26)$$

where  $m_r$  represents the highest similarity value in row  $k$  of the similarity matrix, and  $i^*$  denotes the column of the corresponding infrared cluster. To ensure bidirectional consistency, a reverse selection step is performed. Specifically, for the cluster linked by  $i^*$ , we identify the maximum similarity value in the corresponding column:

$$m_c = \max_k \mathbf{S}_{ki^*}, \quad (27)$$

where  $m_c$  denotes the maximum similarity value in column  $i^*$  of the similarity matrix.

The final label assignment for the  $k$ -th visible instance is determined based on a comparison between the two maximum similarity values  $m_r$  and  $m_c$ . We define that if  $m_r > \gamma \cdot m_c$ , the visible instance is assigned the label corresponding to column  $i^*$ . Here, the parameter  $\gamma$  serves as a critical threshold that controls the selection criteria, ensuring that only sufficiently confident matches are retained. The set of modality-unified pseudo-labels is defined as  $\hat{\mathcal{Y}}^v = \{\hat{\mathbf{y}}_i^v\}_{i=1}^{M_v}$ .

For each infrared instance feature, we perform the same operations for filtering and matching, thus ensuring a bidirectional approach. Similarly, we can obtain the set of modality-unified pseudo-labels  $\hat{\mathcal{Y}}^r = \{\hat{\mathbf{y}}_i^r\}_{i=1}^{M_r}$ . This mechanism effectively addresses model stagnation by alternately transferring labels between visible and infrared instances through multi-perspective transformations. To ensure the reliability of these transfers, this mechanism filters out unreliable pseudo-label matches, thereby enhancing the quality and robustness.

## IV. EXPERIMENT

## A. Datasets and Evaluation Protocol

a) *Datasets:* We evaluate the proposed HIL framework on two widely-used visible-infrared person ReID datasets, namely SYSU-MM01 [58] and RegDB [59]. The SYSU-MM01 dataset is collected by 6 different cameras (4 RGB cameras and 2 IR cameras), including 287,628 RGB images and 15,792 IR images of 491 identities. The RegDB dataset is captured by a dual-camera system that aligns visible and infrared images. It includes 412 identities, and each has 10 infrared images and 10 visible images.

b) *Evaluation Protocol:* Cumulative Matching Characteristics (CMC), mean Average Precision (mAP), and mean Inverse Negative Penalty (mINP) [60] are adopted as the evaluation metrics. For the RegDB dataset, we randomly select 206 identities for training and use the remaining 206 identities for testing, evaluating the method in Visible-to-Infrared mode and Infrared-to-Visible mode. We calculate the average result obtained from 10 random splits of the training set and testing set. For the SYSU-MM01 dataset, the training set contains 22,258 visible images and 11,909 infrared images of 395 identities. The remaining 96 identities are adopted for testing, containing 3,803 infrared images for the query set and 301 randomly selected visible images for the gallery set, evaluating the method in All Search and Indoor Search modes. We also calculate the average result obtained from 10 random gallery set selections.

## B. Implementation Details

The proposed HIL framework is implemented using PyTorch, with the feature extractor from TransReID [61] serving as the backbone network and augmented with dual contrastive learning [5]. Coarse-grained pseudo-labels are generated at the start of each training epoch using DBSCAN [17], followed by the secondary clustering for the generation of fine-grained pseudo-labels using K-means [38]. The learning rate is initialized to 0.000035 and reduced by 0.1 every 20 epochs.



TABLE I

COMPARISON WITH THE STATE-OF-THE-ART METHODS ON SYSU-MM01 AND REGDB. "GUR\*" DENOTES GUR WITHOUT CAMERA LABELS. SINCE OUR METHOD DOES NOT REQUIRE ANY CAMERA LABEL INFORMATION, WE DO NOT REPORT THE RESULTS OF GUR WITH CAMERA LABELS FOR FAIR COMPARISON.

Methods	Venue	SYSU-MM01						RegDB					
		All Search			Indoor Search			Visible-to-Infrared			Infrared-to-Visible		
		R1	mAP	mINP	R1	mAP	mINP	R1	mAP	mINP	R1	mAP	mINP
Supervised VI-ReID methods													
SPOT [39]	TIP-22	65.34	62.25	48.86	69.42	74.63	70.48	80.35	72.46	56.19	79.37	72.26	56.06
DART [40]	CVPR-22	68.72	66.29	-	72.52	78.17	-	83.60	75.67	-	81.97	73.78	-
FMCNet [21]	CVPR-22	66.34	62.51	-	68.15	74.09	-	89.12	84.43	-	88.38	83.86	-
MAUM [41]	CVPR-22	71.68	68.79	-	76.97	81.94	-	87.87	85.09	-	86.95	84.34	-
TransVI [22]	TCSVT-23	71.36	68.63	-	77.40	81.31	-	96.66	91.22	-	96.30	91.21	-
DEEN [42]	CVPR-23	74.70	71.80	-	80.30	83.30	-	91.10	85.10	-	89.50	83.40	-
CAL [43]	ICCV-23	74.66	71.73	-	79.69	83.68	-	94.51	88.67	-	93.64	87.61	-
SAAI [24]	ICCV-23	75.90	77.03	-	83.20	88.01	-	91.07	91.45	-	92.09	92.01	-
PMCM [44]	IJCAI-23	75.54	71.16	-	81.52	84.33	-	93.09	89.57	-	91.44	87.15	-
STAR [45]	TMM-23	76.07	72.73	-	83.47	85.76	-	94.09	88.75	-	93.30	88.20	-
PartMix [46]	CVPR-23	77.78	74.62	-	81.52	84.38	-	84.93	82.52	-	85.66	82.27	-
SEFL [23]	CVPR-23	77.12	72.33	-	82.07	82.95	-	95.35	89.98	-	97.57	91.41	-
IDKL [25]	CVPR-24	81.42	79.85	-	87.14	89.37	-	94.72	90.19	-	94.22	90.43	-
HOS-Net [47]	AAAI-24	75.60	74.20	-	84.20	86.70	-	94.70	90.40	-	93.30	89.20	-
DMA [48]	TIFS-24	74.57	70.41	56.50	82.85	85.10	-	93.30	88.34	-	91.50	86.80	-
Semi-supervised VI-ReID methods													
OTLA [4]	ECCV-22	48.20	43.90	-	47.40	56.80	-	49.90	41.80	-	49.60	42.80	-
TAA [9]	TIP-23	48.77	42.43	25.37	50.12	56.02	49.96	62.23	56.00	41.51	63.79	56.53	38.99
DPIS [10]	ICCV-23	58.40	55.60	-	63.00	70.00	-	62.30	53.20	-	61.50	52.70	-
Unsupervised VI-ReID methods													
ADCA [5]	MM-22	45.51	42.73	28.29	50.60	59.11	55.17	67.20	64.05	52.67	68.48	63.81	49.62
CHCR [12]	TCSVT-23	47.72	45.34	-	50.12	42.17	-	68.18	63.75	-	69.08	63.95	-
MBCCM [13]	MM-23	53.14	48.16	32.41	55.21	61.98	57.13	83.79	77.87	65.04	82.82	76.74	61.73
DOTLA [49]	MM-23	50.36	47.36	32.40	53.47	61.73	57.35	85.63	76.71	61.58	82.91	74.97	58.60
CCLNet [50]	MM-23	54.03	50.19	-	56.68	65.12	-	69.94	65.53	-	70.17	66.66	-
PGM [34]	CVPR-23	57.27	51.78	34.96	56.23	62.74	58.13	69.48	65.41	-	69.85	65.17	-
GUR* [51]	ICCV-23	60.95	56.99	41.85	64.22	69.49	64.81	73.91	70.23	58.88	75.00	69.94	56.21
MIMR [52]	KBS-24	46.56	45.88	-	52.26	60.93	-	68.76	64.33	-	68.76	63.83	-
SCA-RCP [53]	TKDE-24	51.41	48.52	33.56	56.77	64.19	59.25	85.59	78.12	-	82.41	75.73	-
BCGM [54]	MM-24	61.70	56.10	38.70	60.90	66.50	62.30	86.80	81.70	68.60	86.70	82.30	71.10
RPNR [37]	MM-24	65.20	60.00	-	68.90	74.40	-	90.90	84.70	-	90.10	83.20	-
MMM [16]	ECCV-24	61.60	57.90	-	64.40	70.40	-	89.70	80.50	-	85.80	77.00	-
PCLHD [36]	NIPS-24	64.40	58.70	-	69.50	74.40	-	84.30	80.70	-	82.70	78.40	-
MULT [14]	IJCV-24	65.03	58.62	42.77	65.35	71.24	66.60	91.50	83.73	69.13	89.08	80.88	64.03
SDCL [35]	CVPR-24	64.49	63.24	51.06	71.37	76.90	73.50	86.91	78.92	62.83	85.76	77.25	59.57
IMSL [55]	TCSVT-24	57.96	53.93	-	58.30	64.31	-	70.08	66.30	-	70.67	66.35	-
PCAL [56]	TIFS-25	57.94	52.85	36.90	60.07	66.73	62.09	86.43	82.51	72.33	86.21	81.23	68.71
SALCR [57]	IJCV-25	64.44	60.44	45.19	67.17	72.88	68.73	90.58	83.87	70.76	88.69	82.66	66.89
Ours	-	66.30	64.95	52.62	71.81	77.52	74.25	92.82	86.61	73.69	92.24	85.43	70.87

The images are resized to  $288 \times 144$  before being fed into the network. Each training batch consists of 8 pseudo-identities, with 16 instances sampled per pseudo-identity for each modality.  $K$  and  $\gamma$  are set to 9 and 0.5 for SYSU-MM01 and 2 and 0.8 for RegDB. The temperature factor  $\tau$  is set to 0.05 and the momentum updating factor  $\mu$  is set to 0.1. Before standard training, we train the SDCL [35] framework for the initial 30 epochs. Subsequently, we continue training our framework for 30 epochs based on this pre-trained model, using SGD as the optimizer.

### C. Comparison with State-of-the-art Methods

To comprehensively evaluate our method in the VI-ReID task, we compare our method with 15 supervised methods, 3 semi-supervised methods, and 18 unsupervised methods. The comparison results on the SYSU-MM01 and RegDB datasets are reported in Table I.

a) *Comparison with Supervised VI-ReID Methods:* Our proposed HIL method demonstrates competitive performance compared to the supervised method FMCNet [21] on the SYSU-MM01 dataset and achieves performance close to state-of-the-art supervised methods on RegDB. Although supervised methods benefit from precise manual annotations, our approach achieves comparable results, showcasing its robustness without relying on labeled data.

b) *Comparison with Semi-Supervised VI-ReID Methods:* Semi-supervised VI-ReID methods aim to minimize labeling costs by utilizing partially labeled data while maintaining performance. In the realm of existing semi-supervised methods, our proposed approach surpasses the three main methods reported in the literature, achieving state-of-the-art performance without relying on any manual annotations. By eliminating the need for labeled data, our method demonstrates its ability to significantly reduce the dependency on human annotation

TABLE II  
ABLATION STUDY ON THE INDIVIDUAL COMPONENTS OF OUR METHOD ON SYSU-MM01 AND REGDB.

ID	Components			SYSU-MM01						RegDB					
	Baseline	MCCL	BRST	All Search			Indoor Search			Visible-to-Infrared			Infrared-to-Visible		
				R1	mAP	mINP	R1	mAP	mINP	R1	mAP	mINP	R1	mAP	mINP
1	✓			54.88	52.62	38.06	59.02	66.14	61.85	84.11	70.44	48.45	81.66	67.94	45.31
2	✓	✓		55.83	53.76	39.36	60.67	67.62	63.39	87.16	75.94	56.36	85.05	73.50	52.92
3	✓		✓	65.53	63.60	50.67	71.06	76.54	73.08	92.23	85.34	70.88	91.52	84.12	68.07
4	✓	✓	✓	<b>66.30</b>	<b>64.95</b>	<b>52.62</b>	<b>71.81</b>	<b>77.52</b>	<b>74.25</b>	<b>92.82</b>	<b>86.61</b>	<b>73.69</b>	<b>92.24</b>	<b>85.43</b>	<b>70.87</b>

TABLE III  
COMPARISON OF CROSS-MODAL LABEL ASSOCIATION METHODS FOR UNSUPERVISED VI-ReID ON SYSU-MM01 AND REGDB.

ID	Methods	SYSU-MM01						RegDB					
		All Search			Indoor Search			Visible-to-Infrared			Infrared-to-Visible		
		R1	mAP	mINP	R1	mAP	mINP	R1	mAP	mINP	R1	mAP	mINP
1	Baseline	54.88	52.62	38.06	59.02	66.14	61.85	84.11	70.44	48.45	81.66	67.94	45.31
2	Baseline + OPTM [37]	64.31	62.88	50.34	70.27	76.00	72.51	91.33	84.31	70.04	90.76	83.32	67.73
3	Baseline + CRA [35]	65.04	63.35	50.51	70.73	76.36	72.93	91.82	84.80	70.63	91.04	83.75	67.89
4	Baseline + BRST (Ours)	<b>65.53</b>	<b>63.60</b>	<b>50.67</b>	<b>71.06</b>	<b>76.54</b>	<b>73.08</b>	<b>92.23</b>	<b>85.34</b>	<b>70.88</b>	<b>91.52</b>	<b>84.12</b>	<b>68.07</b>

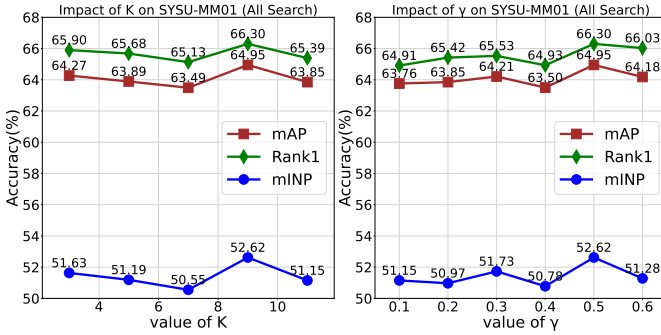


Fig. 3. Hyperparameter analysis of  $K$  and  $\gamma$  on SYSU-MM01.

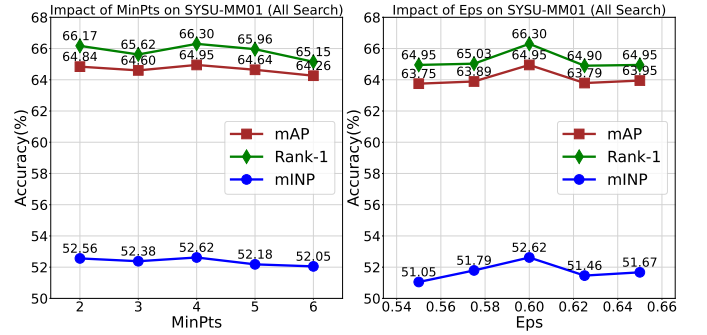


Fig. 4. Hyperparameter analysis of MinPts and Eps on SYSU-MM01.

efforts, while still achieving competitive and robust results across challenging cross-modal VI-ReID tasks.

#### c) Comparison with Unsupervised VI-ReID Methods:

Our method outperforms state-of-the-art unsupervised VI-ReID approaches by a considerable margin. Specifically, it achieves 64.95% mAP and 66.30% Rank-1 on SYSU-MM01 (All Search) and 86.61% mAP and 92.82% Rank-1 on RegDB (Visible-to-Infrared). Compared to the best-performing unsupervised method, SDCL [35], our method exceeds it by 1.81% on SYSU-MM01 (All Search) and 5.91% on RegDB (Visible-to-Infrared) in Rank-1, and by 1.71% on SYSU-MM01 (All Search) and 7.69% on RegDB (Visible-to-Infrared) in mAP. These results underscore the effectiveness of our approach in building robust cross-modal relationships and addressing noise in cross-modal label correspondences.

#### D. Ablation Study

In this subsection, we conduct ablation experiments to validate the effectiveness of each component in our method. The results are reported in Table II.

a) *Baseline*: We denotes the augmented dual-contrastive learning framework [5] with a dual-path transformer architecture. The network optimization is conducted using only the identity loss  $\mathcal{L}_{id}$  and the neighbor loss  $\mathcal{L}_{neighbor}$  during training.

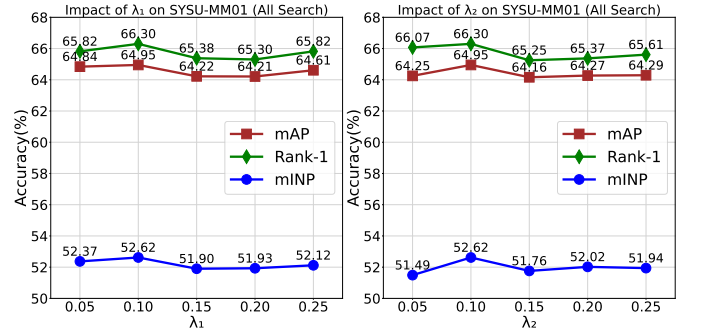


Fig. 5. Hyperparameter analysis of  $\lambda_1$  and  $\lambda_2$  on SYSU-MM01.

b) *Effectiveness of MCCL*: The effectiveness of MCCL is validated through comprehensive comparisons. Specifically, when comparing row 1 and row 2, we observe significant performance improvements on both SYSU-MM01 and RegDB benchmarks in different settings. On SYSU-MM01, MCCL achieves a notable gain of +1.14%/+0.95% in mAP/Rank-1, while on RegDB, it leads to an even more substantial improvement of +5.50%/+3.05% in mAP/Rank-1 compared to the baseline. Similarly, when comparing row 3 and row 4, the inclusion of MCCL yields further enhancements. For SYSU-MM01 (All Search), MCCL provides an addi-



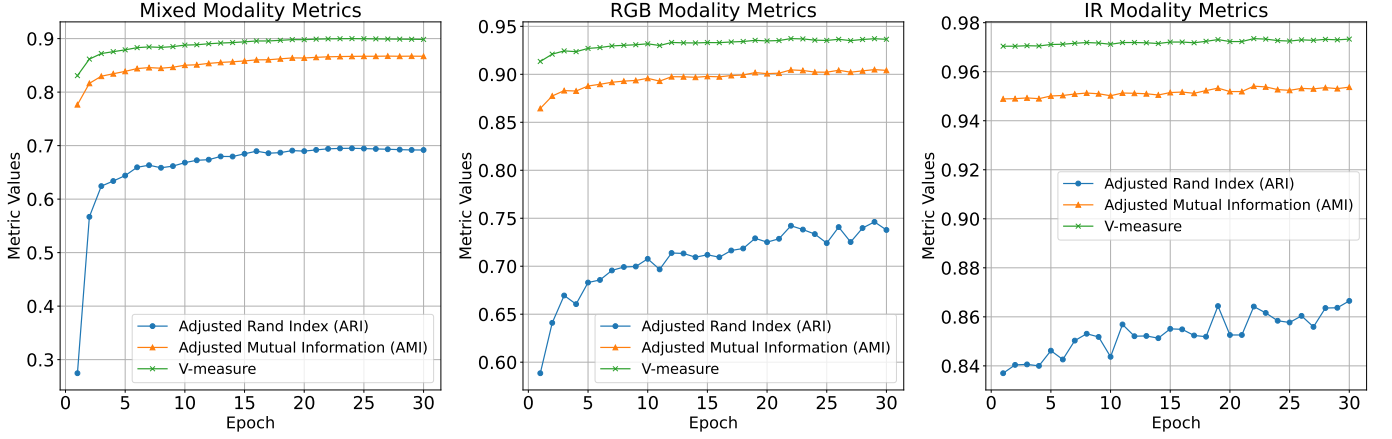


Fig. 6. The variation of clustering quality produced by DBSCAN on the SYSU-MM01 dataset with epoch. The RGB modality represents the clustering quality of pseudo-labels generated using only visible light image modality data. The IR modality indicates the clustering quality of pseudo-labels generated using only infrared image modality data. The mixed modality refers to the overall pseudo-label clustering quality formed by fusing visible light and infrared image modality data.

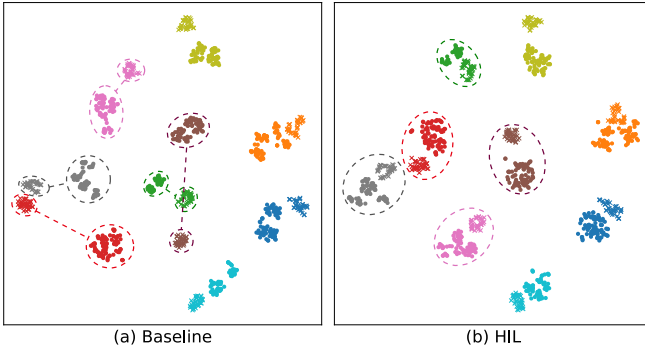


Fig. 7. T-SNE visualization of features for randomly selected identities. Specifically, colors indicate identities, with circular markers denoting the visible modality and cross markers indicating the infrared modality.

tional boost of  $+1.35\%/+0.77\%$  in mAP/Rank-1. On RegDB (Visible-to-Infrared), the improvements are  $+1.27\%/+0.59\%$  in mAP/Rank-1. These results clearly demonstrate that the proposed MCCL effectively leverages hierarchical identity information, thereby reducing cross-modal discrepancies and improving the quality of feature representations across modalities.

*c) Effectiveness of BRST:* The effectiveness of BRST is demonstrated when comparing row 1 and row 3, with improvements of  $+10.98\%/+10.65\%$  mAP/Rank-1 on SYSU-MM01 (All Search) and  $+14.90\%/+8.12\%$  mAP/Rank-1 on RegDB (Visible-to-Infrared) compared to the baseline. Additionally, when comparing row 2 and row 4, we observe improvements of  $+11.19\%/+10.47\%$  mAP/Rank-1 on SYSU-MM01 (All Search) and  $+10.67\%/+5.66\%$  mAP/Rank-1 on RegDB (Visible-to-Infrared) by using BRST. To further evaluate the effectiveness of BRST, we compare BRST with other cross-modal label association methods in Table III. The results demonstrate that our BRST provides higher-quality associations than other methods, ensuring robust cross-modal correspondences by filtering unreliable pseudo-label matches through a reverse selection process.

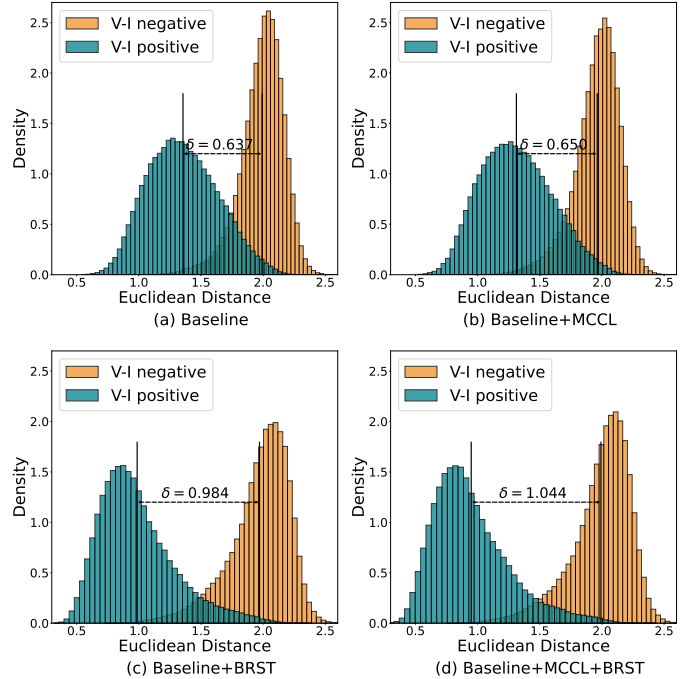


Fig. 8. Visualization of Euclidean Distance distribution of randomly selected visual-infrared positive and negative pairs with various combinations.  $\delta$  represents the average margin between the distances of positive pairs and negative pairs. As  $\delta$  increases, the separability of positive and negative pairs in the cross-modal setting becomes stronger.

## E. Further Analysis

*a) Hyperparameter Analysis:* There are two critical hyperparameters in our method,  $K$  and  $\gamma$ , whose impacts on performance are quantitatively evaluated across a range of values, as shown in Figure 3. The left panel of the figure illustrates the effect of  $K$ , where mAP, Rank-1 accuracy, and mINP all reach their highest values when  $K$  is set to 9. This setting achieves an optimal balance between cross-modal alignment and noise suppression. Similarly, the right panel highlights the influence of  $\gamma$ , showing that the best



Fig. 9. Visualization of the ranking lists on the SYSU-MM01 dataset. The persons who are different from the query persons are marked with red boxes, while those who are the same as the query are marked with green boxes.

performance is achieved when  $\gamma$  is set to 0.5. At this value, the method maintains robust bidirectional selection while avoiding excessively strict or overly lenient matching criteria. These optimal hyperparameter values maximize accuracy across all evaluation metrics, underscoring their significance in enhancing the overall performance of our approach.

In DBSCAN, the two key hyperparameters MinPts and Eps and their impacts on performance have been quantitatively assessed across a range of values, as illustrated in Figure 4. MinPts is the minimum number of points to form a dense cluster, and Eps is the maximum neighborhood distance. The left panel demonstrates that when MinPts is set to 4, mAP, Rank-1 accuracy, and mINP all reach their highest values, where a smaller MinPts enhances sensitivity to local feature clusters but may introduce noise if too small, while larger values risk losing fine-grained features and reducing accuracy. Similarly, the right panel highlights that the best performance occurs when Eps is set to 0.6, as an appropriate Eps balances the model’s ability to capture valid features and filter interference. These results highlight the critical role of hyperparameter tuning in optimizing clustering performance for DBSCAN.

In terms of loss, our method includes two crucial hyperparameters,  $\lambda_1$  and  $\lambda_2$ , whose impacts on performance are quantitatively evaluated across a range of values, as shown in Figure 5. The left panel of the figure illustrates the effect of  $\lambda_1$ , where mAP, Rank-1 accuracy, and mINP all reach their peak values when  $\lambda_1$  is set to 0.10. Similarly, the right panel demonstrates that the highest performance across all three metrics is achieved when  $\lambda_2$  is set to 0.1. When both  $\lambda_1$  and

$\lambda_2$  are set to 0.1, these hyperparameters effectively balance the contributions of  $\mathcal{L}_{id}$ ,  $\mathcal{L}_{neighbor}$ , and  $\mathcal{L}_{mccl}$ , thereby enhancing both the discriminative power and generalization ability of the model.

*b) Clustering Quality:* The overall performance is closely linked to the quality of clustering. As depicted in Fig.6, the initial clusters produced by DBSCAN on the SYSU-MM01 dataset exhibit suboptimal performance. However, clustering quality improves significantly throughout iterative training. Concretely, metrics including Adjusted Rand Index (ARI) [62], Adjusted Mutual Information (AMI) [63], and V-measure [64] for all modalities all exhibit an upward trend as the number of training epochs increases. Specifically, MCCL enhances intra-modal clustering by refining feature representations using fine-grained prototypes. At the same time, BRST improves cross-modal clustering by filtering out unreliable pseudo-label matches during label association, leading to more accurate and stable cross-modal alignment.

*c) Complexity of Implementation:* Assume that the number of instances is  $N$ , the number of clusters is  $M$ , the number of subcenters in secondary clustering is  $K$ , and the dimensionality of features is  $D$ . The spatial complexity of the BRST module is  $O(N + M)$  and its temporal complexity is  $O(NMD)$ . For MCCL, the spatial complexity is  $O(NMK)$ , while the temporal complexity is  $O(NMKD)$ . The model contains 98.61 million parameters and requires 32.78 GFLOPs per image for inference.

*d) Visualization:* To demonstrate the effectiveness of our method in learning modality-invariant features, we randomly select 9 identities from the SYSU-MM01 dataset and visu-

alize their feature embeddings using t-SNE [65]. As shown in Figure 7, our approach generates more compact feature distributions for the same identities within each modality compared to the baseline. This indicates that our method effectively reduces intra-modality variations. Furthermore, the feature embeddings of the same identities across different modalities are significantly closer in proximity, demonstrating improved alignment between visible and infrared features.

To further evaluate the performance of our method, we analyze the distributions of Euclidean distances between randomly selected positive and negative visible-infrared pairs. As shown in Figure 8, the variations in  $\delta$  demonstrate the effectiveness of our approach in improving cross-modal separability. Specifically, positive cross-modal pairs exhibit a higher degree of convergence, as indicated by their smaller distances, while negative pairs show clear divergence, with larger distances. This suggests that our method with proposed modules effectively enhances the separability between positive and negative pairs in the cross-modal setting.

In addition, we visualize some of the ranking lists on the SYSU-MM01 dataset, as shown in Figure 9, and compare the performance of our method with the baseline. The results further demonstrate the effectiveness of our approach in generating high-quality cross-modality pseudo-labels, leading to more accurate ranking outcomes.

## V. CONCLUSION

We propose a novel Hierarchical Identity Learning (HIL) framework for visible-infrared person re-identification tasks. HIL leverages Multi-Center Contrastive Learning (MCCL) to enhance intra-modal clustering and minimize cross-modal discrepancies by refining representations through contrastive learning. Additionally, the Bidirectional Reverse Selection Transmission (BRST) mechanism improves cross-modal matching by performing bidirectional pseudo-label matching and filtering unreliable pseudo-label matches. Extensive experiments demonstrate that our approach achieves superior performance in various settings. Our future work could focus on adaptive thresholding and dynamic refinement of pseudo-labels within the BRST mechanism, which could be investigated to improve the accuracy of label assignments.

## REFERENCES

- [1] Z. Feng, J. Lai, and X. Xie, "Learning modality-specific representations for visible-infrared person re-identification," *IEEE Transactions on Image Processing*, vol. 29, pp. 579–590, 2019.
- [2] M. Ye, X. Lan, Q. Leng, and J. Shen, "Cross-modality person re-identification via modality-aware collaborative ensemble learning," *IEEE Transactions on Image Processing*, vol. 29, pp. 9387–9399, 2020.
- [3] M. Ye, W. Ruan, B. Du, and M. Z. Shou, "Channel augmented joint learning for visible-infrared recognition," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 13 567–13 576.
- [4] J. Wang, Z. Zhang, M. Chen, Y. Zhang, C. Wang, B. Sheng, Y. Qu, and Y. Xie, "Optimal transport for label-efficient visible-infrared person re-identification," in *European Conference on Computer Vision*. Springer, 2022, pp. 93–109.
- [5] B. Yang, M. Ye, J. Chen, and Z. Wu, "Augmented dual-contrastive aggregation learning for unsupervised visible-infrared person re-identification," in *ACM MM*, 2022, p. 2843–2851.
- [6] X. Zheng, X. Chen, and X. Lu, "Visible-infrared person re-identification via partially interactive collaboration," *IEEE Transactions on Image Processing*, vol. 31, pp. 6951–6963, 2022.
- [7] H. Wang, J. Shen, Y. Liu, Y. Gao, and E. Gavves, "Nformer: Robust person re-identification with neighbor transformer," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 7297–7307.
- [8] Y. Cho, W. J. Kim, S. Hong, and S.-E. Yoon, "Part-based pseudo label refinement for unsupervised person re-identification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 7308–7318.
- [9] B. Yang, J. Chen, X. Ma, and M. Ye, "Translation, association and augmentation: Learning cross-modality re-identification from single-modality annotation," *IEEE Transactions on Image Processing*, 2023.
- [10] J. Shi, Y. Zhang, X. Yin, Y. Xie, Z. Zhang, J. Fan, Z. Shi, and Y. Qu, "Dual pseudo-labels interactive self-training for semi-supervised visible-infrared person re-identification," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 11 218–11 228.
- [11] W. Liang, G. Wang, J. Lai, and X. Xie, "Homogeneous-to-heterogeneous: Unsupervised learning for rgb-infrared person re-identification," *IEEE Transactions on Image Processing*, vol. 30, pp. 6392–6407, 2021.
- [12] Z. Pang, C. Wang, L. Zhao, Y. Liu, and G. Sharma, "Cross-modality hierarchical clustering and refinement for unsupervised visible-infrared person re-identification," *IEEE Transactions on Circuits and Systems for Video Technology*, pp. 1–1, 2023.
- [13] D. Cheng, L. He, N. Wang, S. Zhang, Z. Wang, and X. Gao, "Efficient bilateral cross-modality cluster matching for unsupervised visible-infrared person re-id," in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 1325–1333.
- [14] L. He, D. Cheng, N. Wang, and X. Gao, "Exploring homogeneous and heterogeneous consistent label associations for unsupervised visible-infrared person re-id," *arXiv preprint arXiv:2402.00672*, 2024.
- [15] Y. Li, Y. Qin, Y. Sun, D. Peng, X. Peng, and P. Hu, "Romo: Robust unsupervised multimodal learning with noisy pseudo labels," *IEEE Transactions on Image Processing*, 2024.
- [16] J. Shi, X. Yin, Y. Chen, Y. Zhang, Z. Zhang, Y. Xie, and Y. Qu, "Multi-memory matching for unsupervised visible-infrared person re-identification," in *European Conference on Computer Vision*. Springer, 2025, pp. 456–474.
- [17] M. Ester, H.-P. Kriegel, J. Sander, X. Xu *et al.*, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *kdd*, vol. 96, 1996, pp. 226–231.
- [18] Q. Wu, P. Dai, J. Chen, C.-W. Lin, Y. Wu, F. Huang, B. Zhong, and R. Ji, "Discover cross-modality nuances for visible-infrared person re-identification," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 4330–4339.
- [19] M. Ye, W. Ruan, B. Du, and M. Z. Shou, "Channel augmented joint learning for visible-infrared recognition," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 13 567–13 576.
- [20] X. Li, Y. Lu, B. Liu, Y. Liu, G. Yin, Q. Chu, J. Huang, F. Zhu, R. Zhao, and N. Yu, "Counterfactual intervention feature transfer for visible-infrared person re-identification," in *European conference on computer vision*. Springer, 2022, pp. 381–398.
- [21] Q. Zhang, C. Lai, J. Liu, N. Huang, and J. Han, "Fmcnet: Feature-level modality compensation for visible-infrared person re-identification," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 7349–7358.
- [22] Z. Chai, Y. Ling, Z. Luo, D. Lin, M. Jiang, and S. Li, "Dual-stream transformer with distribution alignment for visible-infrared person re-identification," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 33, no. 11, pp. 6764–6776, 2023.
- [23] J. Feng, A. Wu, and W.-S. Zheng, "Shape-erased feature learning for visible-infrared person re-identification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 22 752–22 761.
- [24] X. Fang, Y. Yang, and Y. Fu, "Visible-infrared person re-identification via semantic alignment and affinity inference," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 11 270–11 279.
- [25] K. Ren and L. Zhang, "Implicit discriminative knowledge learning for visible-infrared person re-identification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 393–402.

- [26] Y. Ge, D. Chen, and H. Li, "Mutual mean-teaching: Pseudo label refinery for unsupervised domain adaptation on person re-identification," *arXiv preprint arXiv:2001.01526*, 2020.
- [27] Y. Ge, F. Zhu, D. Chen, R. Zhao *et al.*, "Self-paced contrastive learning with hybrid memory for domain adaptive object re-id," *Advances in neural information processing systems*, vol. 33, pp. 11 309–11 321, 2020.
- [28] H. Chen, B. Lagadec, and F. Bremond, "Ice: Inter-instance contrastive encoding for unsupervised person re-identification," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 14 960–14 969.
- [29] X. Zhang, Y. Ge, Y. Qiao, and H. Li, "Refining pseudo labels with clustering consensus over generations for unsupervised object re-identification," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 3436–3445.
- [30] S. Xuan and S. Zhang, "Intra-inter camera similarity for unsupervised person re-identification," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 11 926–11 935.
- [31] X. Zhang, D. Li, Z. Wang, J. Wang, E. Ding, J. Q. Shi, Z. Zhang, and J. Wang, "Implicit sample extension for unsupervised person re-identification," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 7369–7378.
- [32] Z. Dai, G. Wang, W. Yuan, S. Zhu, and P. Tan, "Cluster contrast for unsupervised person re-identification," in *Proceedings of the Asian Conference on Computer Vision (ACCV)*, December 2022, pp. 1142–1160.
- [33] D. Cheng, X. Huang, N. Wang, L. He, Z. Li, and X. Gao, "Unsupervised visible-infrared person reid by collaborative learning with neighbor-guided label refinement," in *Proceedings of the 31st ACM international conference on multimedia*, 2023, pp. 7085–7093.
- [34] Z. Wu and M. Ye, "Unsupervised visible-infrared person re-identification via progressive graph matching and alternate learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 9548–9558.
- [35] B. Yang, J. Chen, and M. Ye, "Shallow-deep collaborative learning for unsupervised visible-infrared person re-identification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 16 870–16 879.
- [36] J. Shi, X. Yin, Y. Zhang, Y. Xie, Y. Qu *et al.*, "Learning commonality, divergence and variety for unsupervised visible-infrared person re-identification," in *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [37] X. Yin, J. Shi, Y. Zhang, Y. Lu, Z. Zhang, Y. Xie, and Y. Qu, "Robust pseudo-label learning with neighbor relation for unsupervised visible-infrared person re-identification," in *Proceedings of the 32nd ACM International Conference on Multimedia*, 2024, pp. 2242–2251.
- [38] J. MacQueen *et al.*, "Some methods for classification and analysis of multivariate observations," in *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, vol. 1. Oakland, CA, USA, 1967, pp. 281–297.
- [39] C. Chen, M. Ye, M. Qi, J. Wu, J. Jiang, and C.-W. Lin, "Structure-aware positional transformer for visible-infrared person re-identification," *IEEE Transactions on Image Processing*, vol. 31, pp. 2352–2364, 2022.
- [40] M. Yang, Z. Huang, P. Hu, T. Li, J. Lv, and X. 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 (CVPR)*, June 2022, pp. 14 308–14 317.
- [41] J. Liu, Y. Sun, F. Zhu, H. Pei, Y. Yang, and W. Li, "Learning memory-augmented unidirectional metrics for cross-modality person re-identification," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 19 366–19 375.
- [42] Y. Zhang and H. Wang, "Diverse embedding expansion network and low-light cross-modality benchmark for visible-infrared person re-identification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 2153–2162.
- [43] J. Wu, H. Liu, Y. Su, W. Shi, and H. Tang, "Learning concordant attention via target-aware alignment for visible-infrared person re-identification," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 11 122–11 131.
- [44] Z. Qian, Y. Lin, and B. Du, "Visible-infrared person re-identification via patch-mixed cross-modality learning," *Pattern Recognition*, vol. 157, p. 110873, 2025.
- [45] J. Wu, H. Liu, W. Shi, M. Liu, and W. Li, "Style-agnostic representation learning for visible-infrared person re-identification," *IEEE Transactions on Multimedia*, vol. 26, pp. 2263–2275, 2023.
- [46] M. Kim, S. Kim, J. Park, S. Park, and K. Sohn, "Partmix: Regularization strategy to learn part discovery for visible-infrared person re-identification," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 18 621–18 632.
- [47] L. Qiu, S. Chen, Y. Yan, J.-H. Xue, D.-H. Wang, and S. Zhu, "High-order structure based middle-feature learning for visible-infrared person re-identification," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 5, 2024, pp. 4596–4604.
- [48] Z. Cui, J. Zhou, and Y. Peng, "Dma: Dual modality-aware alignment for visible-infrared person re-identification," *IEEE Transactions on Information Forensics and Security*, vol. 19, pp. 2696–2708, 2024.
- [49] D. Cheng, X. Huang, N. Wang, L. He, Z. Li, and X. Gao, "Unsupervised visible-infrared person reid by collaborative learning with neighbor-guided label refinement," 2023.
- [50] Z. Chen, Z. Zhang, X. Tan, Y. Qu, and Y. Xie, "Unveiling the power of clip in unsupervised visible-infrared person re-identification," in *Proceedings of the 31st ACM International Conference on Multimedia*, ser. MM '23. New York, NY, USA: Association for Computing Machinery, 2023, p. 3667–3675. [Online]. Available: <https://doi.org/10.1145/3581783.3612050>
- [51] B. Yang, J. Chen, and M. Ye, "Towards grand unified representation learning for unsupervised visible-infrared person re-identification," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 11 069–11 079.
- [52] Z. Pang, C. Wang, H. Pan, L. Zhao, J. Wang, and M. Guo, "Mimr: Modality-invariance modeling and refinement for unsupervised visible-infrared person re-identification," *Knowledge-Based Systems*, vol. 285, p. 111350, 2024.
- [53] Z. Li, H. Liu, X. Peng, and W. Jiang, "Inter-intra modality knowledge learning and clustering noise alleviation for unsupervised visible-infrared person re-identification," *IEEE Transactions on Knowledge and Data Engineering*, 2024.
- [54] X. Teng, X. Shen, K. Xu, and L. Lan, "Enhancing unsupervised visible-infrared person re-identification with bidirectional-consistency gradual matching," in *Proceedings of the 32nd ACM International Conference on Multimedia*, 2024, pp. 9856–9865.
- [55] Z. Pang, L. Zhao, Y. Liu, G. Sharma, and C. Wang, "Inter-modality similarity learning for unsupervised multi-modality person re-identification," *IEEE Transactions on Circuits and Systems for Video Technology*, 2024.
- [56] Y. Yang, W. Hu, and H. Hu, "Progressive cross-modal association learning for unsupervised visible-infrared person re-identification," *IEEE Transactions on Information Forensics and Security*, 2025.
- [57] D. Cheng, L. He, N. Wang, D. Zhang, and X. Gao, "Semantic-aligned learning with collaborative refinement for unsupervised vi-reid," *International Journal of Computer Vision*, pp. 1–23, 2025.
- [58] A. Wu, W.-S. Zheng, H.-X. Yu, S. Gong, and J. Lai, "Rgb-infrared cross-modality person re-identification," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 5380–5389.
- [59] D. T. Nguyen, H. G. Hong, K. W. Kim, and K. R. Park, "Person recognition system based on a combination of body images from visible light and thermal cameras," *Sensors*, vol. 17, no. 3, p. 605, 2017.
- [60] M. Ye, J. Shen, G. Lin, T. Xiang, L. Shao, and S. C. Hoi, "Deep learning for person re-identification: A survey and outlook," *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 6, pp. 2872–2893, 2021.
- [61] S. He, H. Luo, P. Wang, F. Wang, H. Li, and W. Jiang, "Transreid: Transformer-based object re-identification," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 15 013–15 022.
- [62] L. Hubert and P. Arabie, "Comparing partitions journal of classification 2 193–218," *Google Scholar*, vol. 193, 1985.
- [63] N. X. Vinh, J. Epps, and J. Bailey, "Information theoretic measures for clusterings comparison: is a correction for chance necessary?" in *Proceedings of the 26th annual international conference on machine learning*, 2009, pp. 1073–1080.
- [64] J. B. Hirschberg and A. Rosenberg, "V-measure: a conditional entropy-based external cluster evaluation," 2007.
- [65] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne," *Journal of machine learning research*, vol. 9, no. 11, 2008.





**Haonan Shi** is currently pursuing a B.Sc. degree in Intelligent Science and Technology at Xidian University, Xi'an, China. His research interests include computer vision, pattern recognition, and machine learning.



**Yubin Wang** received the B.E. degree in data science and big data technology from Tongji University, China, in 2022, where he is currently pursuing the master's degree. His main research interests include prompt learning, multi-modal learning, and person re-identification.



**De Cheng** is an associate professor with School of Telecommunications Engineering, Xidian University, China. He received the B.S. and Ph.D. degrees from Xi'an Jiaotong University, Xi'an, China, in 2011 and 2017, respectively. From 2015 to 2017, he was a visiting scholar in Carnegie Mellon University, Pittsburgh, USA. His research interests include pattern recognition, machine learning, and multimedia analysis.



**Lingfeng He** received the B.Sc. degree from Xidian University, Xi'an, China, in 2023. He is currently pursuing his M.S. degree in Information and Communication Engineering in Xidian University. His research interests in person re-identification and unsupervised learning.



**Nannan Wang** (M'16) received the B.Sc. degree in information and computation science from the Xi'an University of Posts and Telecommunications in 2009 and the Ph.D. degree in information and telecommunications engineering from Xidian University in 2015. From September 2011 to September 2013, he was a Visiting Ph.D. Student with the University of Technology, Sydney, NSW, Australia. He is currently a Professor with the State Key Laboratory of Integrated Services Networks, Xidian University. He has published over 100 articles in

refereed journals and proceedings, including IEEE T-PAMI, IJCV, CVPR, ICCV etc. His current research interests include computer vision and machine learning.



**Xinbo Gao** (M'02-SM'07) received the B.Eng., M.Sc. and Ph.D. degrees in electronic engineering, signal and information processing from Xidian University, Xi'an, China, in 1994, 1997, and 1999, respectively. From 1997 to 1998, he was a Research Fellow with the Department of Computer Science, Shizuoka University, Shizuoka, Japan. From 2000 to 2001, he was a Post-Doctoral Research Fellow with the Department of Information Engineering, the Chinese University of Hong Kong, Hong Kong. Since 2001, he has been with the School of Elec-

tronic Engineering, Xidian University. He is also a Cheung Kong Professor of the Ministry of Education of China, a Professor of Pattern Recognition and Intelligent System with Xidian University, and a Professor of Computer Science and Technology with the Chongqing University of Posts and Telecommunications, Chongqing, China. He has published 6 books and around 300 technical articles in refereed journals and proceedings. His research interests include image processing, computer vision, multimedia analysis, machine learning, and pattern recognition. Prof. Gao is also a Fellow of the Institute of Engineering and Technology and the Chinese Institute of Electronics. He has served as the general chair/cochair, the program committee chair/co-chair, or a PC member for around 30 major international conferences. He is also on the Editorial Boards of several journals, including Signal Processing (Elsevier) and Neurocomputing (Elsevier).