# Domain Consistency Representation Learning for Lifelong Person Re-Identification

Shiben Liu Ⓘ, Huijie Fan* Ⓘ, Qiang Wang Ⓘ, Weihong Ren Ⓘ, BaojieFan Ⓘ, Yandong Tang Ⓘ

*Abstract*—Lifelong person re-identification (LReID) exhibits a contradictory relationship between intra-domain discrimination and inter-domain gaps when learning from continuous data. Intra-domain discrimination focuses on individual nuances (*i.e.*, clothing type, accessories, *etc.*), while inter-domain gaps emphasize domain consistency. Achieving a trade-off between maximizing intra-domain discrimination and minimizing inter-domain gaps is a crucial challenge for improving LReID performance. Most existing methods strive to reduce inter-domain gaps through knowledge distillation to maintain domain consistency. However, they often ignore intra-domain discrimination. To address this challenge, we propose a novel domain consistency representation learning (DCR) model that explores global and attribute-wise representations as a bridge to balance intra-domain discrimination and inter-domain gaps. At the intra-domain level, we explore the complementary relationship between global and attribute-wise representations to improve discrimination among similar identities. Excessive learning intra-domain discrimination can lead to catastrophic forgetting. We further develop an attribute-oriented anti-forgetting (AF) strategy that explores attribute-wise representations to enhance inter-domain consistency, and propose a knowledge consolidation (KC) strategy to facilitate knowledge transfer. Extensive experiments show that our DCR model achieves superior performance compared to state-of-the-art LReID methods. Our code will be available soon.

*Index Terms*—Lifelong learning, person re-identification, domain consistency representations, attribute and text guided representations.

## I. INTRODUCTION

Person re-identification (ReID) aims to retrieve the same individual across non-overlapping cameras in a large-scale database, and has achieved significant progress using unimodal architectures such as convolutional neural networks (CNN) [3, 4] or vision transformers (ViT) [5–7]. However, when ReID models are applied to continuous datasets collected by dynamic monitoring systems, they exhibit notable
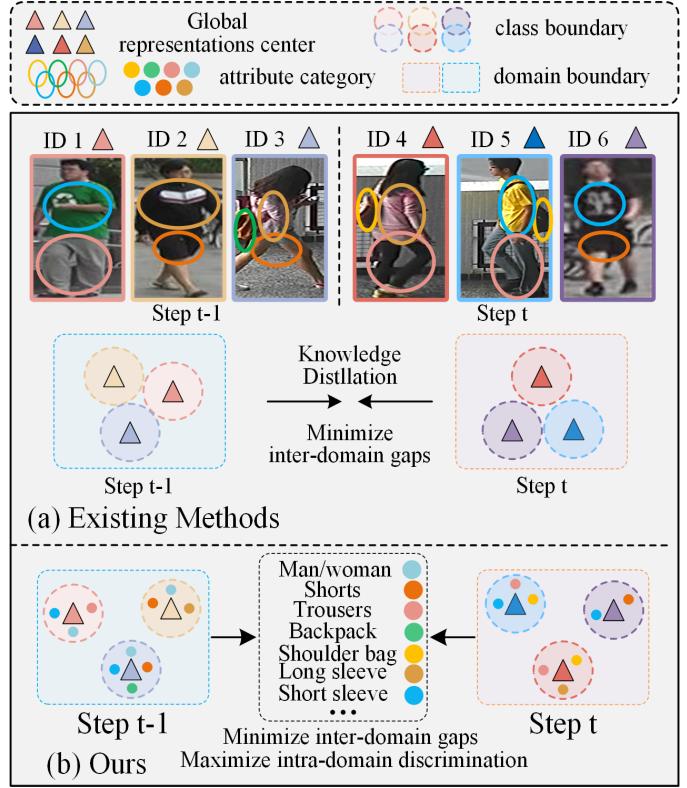
Fig. 1. Comparison between our method and existing methods. (a) Existing methods [1, 2] leverage knowledge distillation to minimize inter-domain gaps but ignore intra-domain discrimination, which limits the LReID model's ability to learn new knowledge. (b) Our method explores attribute-wise representations as a bridge to achieve a trade-off between maximizing intra-domain discrimination and minimizing inter-domain gaps, enhancing the LReID model's anti-forgetting and generalization capabilities.

performance limitations. Thus recent works have focused more on the practical problem of lifelong person re-identification (LReID), which involves learning from streaming data and maintaining strong performance across all data.

At present, lifelong person re-identification (LReID) suffers from the challenge of balancing the anti-forgetting of old knowledge and learning new knowledge. Specifically, there are two main issues to solve this challenge. 1) **Intra-domain discrimination**. Each identity may exhibit subtle nuances of individual information (*i.e.*, clothing type, accessories, haircut, *etc.*) and lead to severe distribution overlapping. Learning discriminative representations of individuals are effective for distinguish identity information. 2) **Inter-domain gaps**. The dataset of each task is collected in different illumination and background, leading to inter-domain gaps. Bridging intra-

domain gaps are significant for mitigating catastrophic forgetting in LReID.

To address these issues, we aim to learn consistency representations that capture individual nuances in intra-domain and inter-domain consistency in LReID, striking a balance between maximizing intra-domain discrimination and minimizing inter-domain gaps. Knowledge distillation-based approaches [2, 8–10] ensure distribution consistency between the previous and current datasets to alleviate catastrophic forgetting. However, these approaches impose strict constraints and ignore intra-domain discrimination, [11–13], as outlined in Figure 1(a). While LReID models significantly improve intra-domain discrimination for the current step, they inevitably damage inter-domain consistency, leading to catastrophic forgetting. Thus, we propose consistency representations as a bridge to achieve a trade-off between maximizing intra-domain discrimination and minimizing inter-domain gaps, improving the anti-forgetting and generalization capabilities of the LReID model, as illustrated in Figure 1(b).

Specifically, we propose a novel domain consistency representation learning (DCR) model, which first explores attribute and text information to enhance LReID performance. Unlike methods [14–16], we develop consistency representations including global and attribute-wise representations to capture individual nuances in intra-domain and inter-domain consistency in LReID. We design an attribute-text generator (ATG) to dynamically generate text-image pairs for each instance, which are then fed into a text-guided aggregation (TGA) network to improve the global representation capability, effectively distinguishing identities in LReID. In addition, the attributes of each instance guide an attribute compensation (ACN) network to generate attribute-wise representations focusing on specific regional information of identities. We consider that attributes can ensure reliability by setting higher thresholds across datasets. Therefore, the generated attribute-wise representations and text for each instance are considered reliable in our model.

In summary, we explore global representations and attribute-wise representations to strike a balance between maximizing identity-discriminative information of intra-domain and minimizing inter-domain gaps. At the intra-domain discrimination level, global representations capture whole-body information, while attribute-wise representations focus on specific regional information. When whole-body appearances or attribute-related information are very similar across identities, we combine global and attribute-wise representations to distinguish among similar identities, maximizing intra-domain discrimination. Perfect learning intra-domain discrimination can lead to catasttophic forgetting. We further develop an attribute-oriented anti-forgetting (AF) strategy that explore attribute-wise representations for bridging inter-domain gaps across continuous datasets. Additionally, knowledge consolidation (KC) is proposed to enable knowledge transfer, improving generalization capabilities. Our contributions are as follows:

- We propose a novel domain consistency representation learning (DCR) model that explores global and attribute-wise representations to capture individual nuances in intra-domain and inter-domain consistency, achieving a trade-off between maximizing intra-domain discrimination and minimizing inter-domain gaps.
- In the intra-domain context, we explore the complementary relationship between global and attribute-wise representations to enhance the discrimination of each identity and adapt to new knowledge.
- In the inter-domain context, we design an attribute-oriented anti-forgetting (AF) and a knowledge consolidation (KC) strategy to minimize inter-domain gaps and facilitate knowledge transfer, improving the LReID model's generalization and anti-forgetting capabilities.

### A. Lifelong Person Re-Identification

Lifelong Person Re-Identification (LReID) faces a formidable challenge, aiming to address the evolving nature of person identification across various scenarios. Some works [12, 17] are proposed to tackle the issue of adapting ReID models over time while retaining knowledge gained from previous distribution. Generally, Pu et al. [17] proposed learnable knowledge graphs that adaptively facilitate the mutual exchange of new and old knowledge, thus achieving knowledge accumulation. Some works [1, 8, 9, 14] aim to extract rich and discriminative representation, mitigating the risk of knowledge forgetting. Pu et al. [8] proposed a meta reconciliation normalization (MRN) for mining meta-knowledge shared across different domains. Meanwhile, ConRFL [9] maintains learnable and consistent features across all seen domains, which improves the discrimination and adaptation ability of the LReID model. In addition, some methods [1, 2, 17] mitigate catastrophic forgetting and enhance model accuracy by using rehearsal-based strategies with images stored from previous tasks. These approaches strive to reduce inter-domain gaps, ensuring distribution consistency across datasets to mitigate catastrophic forgetting. However, this strategy employs strict constraints and ignores intra-domain discrimination, limiting LReID model's performance to learn new knowledge. In this paper, we propose consistency representations as a bridge to achieve a trade-off between maximizing intra-domain discrimination and minimizing inter-domain gaps for improving the anti-forgetting and generalization capabilities of the LReID model.

### B. Vision-Language for Person Re-Identification

The vision-language learning paradigms [18–20] have gained considerable popularity in recent years. Contrastive Language-Image Pre-training (CLIP) [21], establishes a connection between natural language and visual content through the similarity constraint of image-text pair. CLIP has been applied to multiple person re-identification tasks [22, 23], including text-to-image, text-based single-modality, text-based cross modality. Text-to-image methods [24–26] aims to retrieve the target person based on a textual query. Text-based single-modality works [5, 6] leverage text descriptions to generate robust visual features or to integrate the beneficial features of text and images for the person
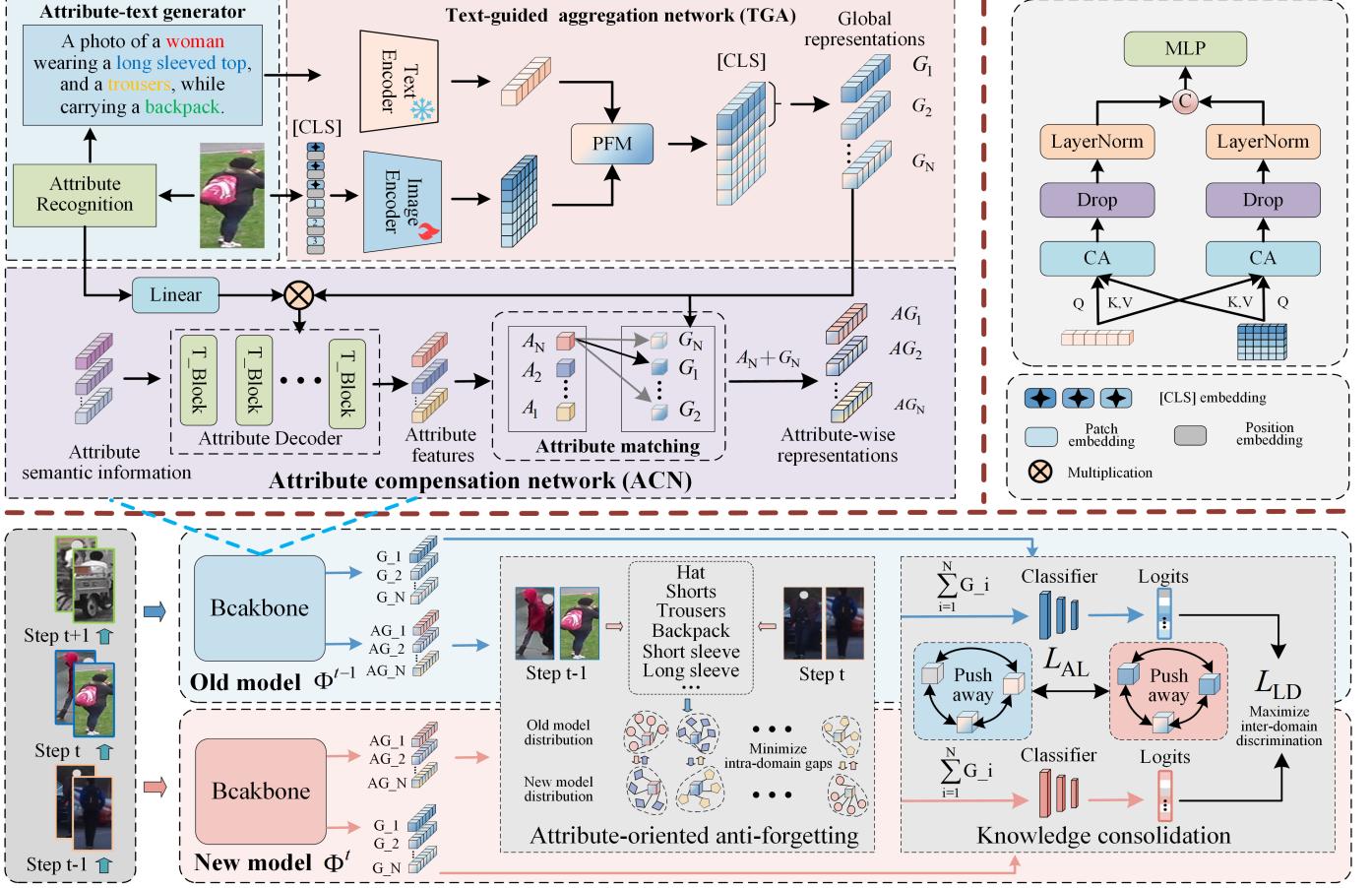
Fig. 2. Overview of the proposed DCR for LReID. First, the attribute-text generator (ATG) dynamically generates text-image pairs for each instance. Then, the text-guided aggregation network (TGA) captures global representations for each identity, while the attribute compensation network (ACN) generates attribute-wise representations. We explore the complementary relationship between global and attribute-wise representations to maximize intra-domain discrimination. Meanwhile, we design attribute-oriented anti-forgetting (AF) and knowledge consolidation (KC) strategies to minimize inter-domain gaps and facilitate knowledge transfer.

category. Text-based cross modality methods [27] employ text descriptions to alleviate visible-infrared modality gaps. Text information generated by prompt learning and text inversion, providing insufficient text descriptions of each identity. In this paper, we dynamically generate text-image pairs from single image to capture fine-grained global representations based on the CLIP model for improving inter-domain discrimination.

### C. Pedestrian Attribute Recognition

Pedestrian attribute recognition aims to assign a set of attributes (Gender, Bag, Short/Long sleeve, and *etc.*) to a visual representation of a pedestrian. Deep learning-based researches [28] automatically learn hierarchical features from raw images, improving recognition accuracy. Multi-task learning approaches [29–31] leverage additional contextual information of across tasks, such as pedestrian detection or pose estimation, to significantly improve attribute recognition. Part-based methods [32, 33] divide the pedestrian image into several parts or regions, providing more accurate attribute localization. At present, the above methods have achieved significant success in improving the accuracy of attribute recognition. We are the

first to explore the application of attributes to LReID from the following two aspects. 1) Attributes are converted into text descriptions for each image to enhance global representation capabilities. 2) Attributes are transformed into attribute-wise representations by specific networks to maximize intra-domain discrimination and minimize intra-domain gaps.

## II. PROPOSED METHOD

### A. Preliminary: Overview of Method

The overview of our DCR model to achieve a trade-off between maximizing intra-domian discrimination and minimizing inter-domain gaps, is depicted as Figure 2. The DCR model learns the old model $\Phi^{t-1}$ and new model $\Phi^t$ from (t-1)-th and t-th steps, where $\Phi^t$ is inherited from $\Phi^{t-1}$. $\Phi^{t-1}$ and $\Phi^t$ with three branches of attribute-text generator (ATG), text-guided aggregation network (TGA) and attribute compensation network (ACN). $\phi^{t-1}$ and $\phi^t$ serve as classifier heads for the old and new models, providing logits of each instance for recognition. Additionally, we define that consecutive $T$ person datasets $D = \{D^t\}_{t=1}^T$ are collected from different environments, and establish a memory buffer $M$ to store a limited number of samples from each previous ReID dataset.

Given an image $x_i^t \in \mathcal{D}^t \cup \mathcal{M}$, we forward it to $\Phi^{t-1}$ and $\Phi^t$ is as follows:

$$G^{t-1}, AG^{t-1} = \Phi^{t-1}(x^i); \quad G^t, AG^t = \Phi^t(x^i). \quad (1)$$

### B. Attribute-Text Generator

Due to the lack of text-image pairs in ReID datasets, we propose an attribute-text generator (ATG) to dynamically generate corresponding text descriptions for each instance. Specifically, we first introduce an attribute recognition model pre-trained on the PA100K dataset [35] to generate attribute categories (*i.e.*, female, backpack, short/long sleeve, and *etc.*), which are then converted into text descriptions for each instance using a specific template. This template adds modifiers (in black font) to each attribute (in a different color font) to create a complete sentence describing an instance, as shown in Figure 2. Although, attributes can vary significantly across datasets, we consider that text descriptions can be made reliable by setting a higher threshold (Confidence threshold=0.80) to ensure classification accuracy of attribute recognition network.

### C. Text-Guided Aggregation Netwrok

We propose a text-guided aggregation network (TGA) to explore global representations for each identity and knowledge transfer, as shown in Figure 2 (TGA). The TGA includes a CLIP model and a parallel fusion module (PFM). Note that the text encoder is frozen in our DCR model.

**Parallel Fusion Module.** By attribute-text generator obtain text-image pairs, we employ CLIP with text encoder $\mathcal{T}(\cdot)$ and image encoder $\mathcal{V}(\cdot)$ to extract text and image embedding, respectively. Unlike CLIP [21], we introduce multiple [CLS] embeddings into the image encoder input sequence to capture multiple global representations from different perspectives. To obtain fine-grained global representations for improving the performance of the LReID model, we propose a parallel fusion module (PFM) to explicitly explore the interactions between image embeddings and text embeddings, as shown in Figure 2 (PFM). Firstly, we leverage text embedding $d^*$ as query and image embedding $[v_1^*, \cdots, v_N^*, v_1, \cdots, v_P]$ as key and value to implement operation with cross-attention, drop, and layer normalization, getting text-wise representations. Similarly, in another fusion branch, image-wise representations are obtained. Finally, image-wise and text-wise representations perform concatenation and MLP operations to obtain global representations $G^t = \{G_i | i = 1, 2, \cdots, N\}$, focusing on whole body information. We force multiple global reprsentatios $G^t$ at the current step to learn more discriminative information by orthogonal loss to minimize the overlapping elements. The orthogonal loss can be formulated as:

$$L_{Ort} = \sum_{i=1}^{N-1} \sum_{j=i+1}^{N} (G_i^t, G_j^t) \quad (2)$$

Then, we utilize the cross-entropy loss $L_{CE}$ and triplet loss $L_{Tri}^g$ [6] to optimize our DCR at the current task.

$$L_{CE} = \frac{1}{K} \sum_{i=1}^{K} y_i \log((\phi^t(G^t))_i) \quad (3)$$

$$L_{Tri}^g = max(d_p^g - d_n^g + m, 0) \quad (4)$$

where $K$ is the number of classes, and $m$ is the margin, $d_p^g$ and $d_n^g$ are the distances from positive samples and negative samples to anchor samples in global representations, respectively. Unlike some methods [2, 11], global representations generated by the text-guided aggregation (TGA) network present two advantages. First, we leverage text descriptions based on the CLIP model to enhance the discrimination capability of global representations, allowing them to better distinguish identities and adapt to new knowledge. Second, global representations facilitate knowledge transfer, improving the model's generalization ability.

### D. Attribute Compensation Network

We force attributes to guide attribute compensation network (ACN) for learning attribute-wise representations. The ACN consists of an attribute decoder and an attribute matching component, as illustrated in Figure 2 (ACN).

**Attribute Decoder.** Enabling attributes to better adapt across datasets, we define multiple learnable attribute semantic information $A^* = \{A_i^* | i = 1, 2, \cdots, N\}$ to learn discriminative information. The attributes undergo a linear layer to increase its dimensions, and then multiplies with the text-image global representation to output $f_{AT}$. Attribute semantic information $A^*$ as queries $Q$, $f_{AT}$ as keys and values are input into attribute decoder, which outputs the attribute features $A = \{A_i | i = 1, 2, \cdots, N\}$. The attribute decoder employs six transformer blocks (T_Block) referenced from [36].

**Attribute Matching.** The attribute features $A = \{A_i | i = 1, 2, \cdots, N\}$ learn multiple discriminative local infromation of individuals. However, it is unclear which attribute features correspond to specific body parts. Thus, we propose an attribute matching (AM) component to associate attribute features and global representations $G = \{G_i | i = 1, 2, \cdots, N\}$. The core objective is to find the most similar global representations$G$ from different perspectives and local attribute features $A$, and then add the them with the highest similarity. Specifically, attribute-wise representations $AG^t = \{AG_i | i = 1, 2, \cdots, N\}$ is formulated as:

$$k = argmax(\frac{<A_i, G>}{|A_i||G|}) \quad (5)$$

$$AG_i = A_i + G_k. \quad (6)$$

We leverage the triplet loss to align attribute-wise representations with identity at the current step, assisting in global representations to distinguish similar identities.

$$L_{Tri}^l = max(d_p - d_n + m, 0) \quad (7)$$

where, $d_p^l$ and $d_n^l$ are the distances from positive samples and negative samples to anchor samples in attribute-wise representations, respectively. In this paper, attribute-wise representations that contain specific information of individuals assist global representations in distinguishing similar identities for maximizing intra-domain discrimination. Meanwhile, attribute-wise representations as a bridge across increasing datasets to minimize inter-domain gaps for better knowledge transfer.
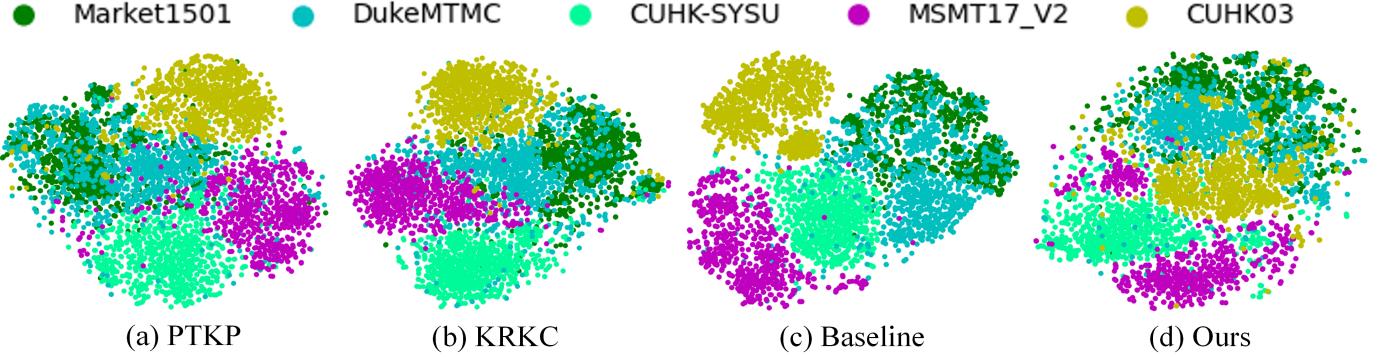
Fig. 3. t-SNE visualization of feature distribution on five seen datasets. Our method better narrows the distribution across datasets for minimizing inter-domain gaps, improving the anti-forgetting and generalization ability of the model.

The knowledge consolidation loss is defined as:

$$L_{KC} = L_{AL} + L_{LD} \tag{11}$$

The total loss function is formulated as:

$$L = L_{CE} + L_{Tri}^g + L_{Tri}^l + L_{Ort} + L_{AF} + L_{KC} \tag{12}$$

## III. EXPERIMENTS

### A. Experiments Setting

**Datasets.** To verify the performance of our method in anti-forgetting and generalization, we evaluate our method on a challenging benchmark consisting of Market1501 [38], CUHK-SYSU [39], DukeMTMC [40], MSMT17_V2 [41] and CUHK03 [42], referred to as seen datasets. Two representative training orders are set up following the protocol described in [17] for training and testing. Further, we employ six datasets including VIPeR [43], GRID [44], CUHK02 [45], Occ_Duke [46], Occ_REID [47], and PRID2011 [48], as unseen dataset. **Implementation Details.** Our text encoder and image encoder are based on a pre-trained CLIP model, while the attribute decoder utilizes a transformer-based architecture[36]. All person images are resized to 256×128. We use Adam [49] for optimization and train each task for 60 epochs. The batch size is set to 128. The learning rate is initialized at $5\times10^{-6}$ and is decreased by a factor of 0.1 every 20 epochs for each task. We employ mean average precision (mAP) and Rank-1 accuracy (R-1) to evaluate the LReID model on each dataset.

### B. Comparison with SOTA Methods

We compare the proposed DCR with SOTA LReID to demonstrate the superiority of our method, including AKA[17], PTKP[1], PatchKD[14], KRKC[2], and ConRFL[9], DKP[15], C2R[16], LSTKC[13]. Experimental results on training order-1 and order-2 are shown in Table I and Table II, respectively. CODA [50] method employs ViT-B/16 as the backbone.
**Compared with LReID methods.** In Table I and Table II, Our DCR significantly outperforms LReID methods, with an seen-avg incremental gain of 10.0% mAP/7.8% R-1, and



Fig. 4. Visualization of intra-domain discrimination on the Market1501 dataset. We randomly select 30 identities. Colors represent different identity information. Our DCR model can cluster images of the same identity more tightly (circle) for minimizing inter-domian discrimination.

9.8% mAP/7.5% R-1 on training order-1 and order-2, respectively. Meanwhile, our DCR effectively alleviate catastrophic forgetting, achieving 6.9% mAp/1.1% R-1, and 5.4% mAP/ 2.2% R-1 improvement on the first dataset (Mrket1501 and DukeMTMC) with different training orders. Compared to CODA, our DCR significantly outpreform performance under the backbone of VIT-B/16. Additionally, our DCR improves the average by 8.1 mAP%/7.5% R-1 and 9.5% mAP/11.0% R-1 on unseen datasets. In contrast, our DCR achieves a trade-off between anti-forgetting and acquiring new information, significantly enhances generalization capabilities.
**Compared with Baseline.** Due to the lack of CLIP-based comparison methods in LReID, we introduce a Baseline model including CLIP model, attribute-text generator and knowledge consolidation strategy. Compared to the Baseline, Our DCR improves the Seen-Avg by 11.4% mAP/10.4% R-1 and by 9.8% mAP/10.2% R-1. These results demonstrate that our proposed attribute-wise representations learning achieves significant performance in balancing maximization intra-domain discrimination and minimization inter-domain gaps in LReID.
**The effectiveness of minimizing inter-domain gaps.** We visualize the feature distribution of PTKP, KRKC, DKP, and our method across five datasets, as shown in Figure 3. DKP shows poor performance in bridging inter-domain gaps, as knowledge prototypes struggle to fit the data distribution. Compared to other methods, our DCR model better narrows the distribution across increasing datasets. Thus, the proposed

Fig. 5. Generalization curves. After each training step, the performance of all unseen domains is evaluated.

TABLE III
ABLATION STUDIES ON THE NUMBER OF GLOBAL AND ATTRIBUTE-WISE REPRESENTATIONS $N$ ON TRAINING ORDER-1.

| Number ($N$) | Seen_Avg | | Unseen_Avg | |
|---|---|---|---|---|
| | mAP | R-1 | mAP | R-1 |
| 2 | 60.2 | 68.7 | 59.4 | 56.5 |
| 3 | **61.8** | **71.9** | **60.8** | **58.3** |
| 4 | 61.2 | 71.6 | 60.3 | 57.5 |

TABLE IV
ABLATION STUDIES OF DIFFERENT COMPONENTS ON TRAINING ORDER-1.

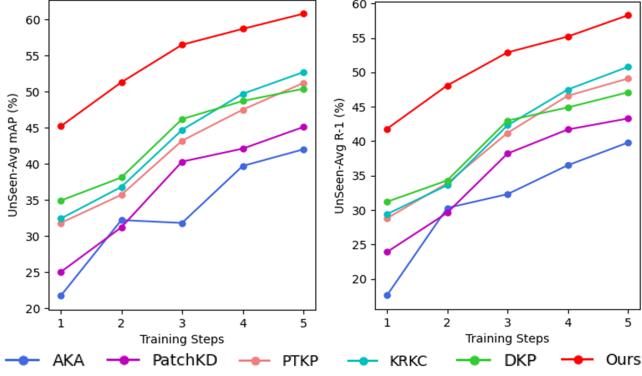| PFM | ACN | AF | KC | Seen_Avg | | Unseen_Avg | |
|---|---|---|---|---|---|---|---|
| | | | | mAP | R-1 | mAP | R-1 |
| | | | | 50.4 | 61.5 | 51.8 | 49.4 |
| √ | | | | 51.7 | 62.1 | 52.5 | 50.3 |
| √ | | | √ | 57.6 | 68.9 | 58.2 | 56.2 |
| √ | √ | √ | | 58.7 | 69.2 | 58.5 | 56.8 |
| √ | √ | √ | √ | **61.8** | **71.9** | **60.8** | **58.3** |

TABLE V
ABLATION OF TRAINING WITH OR WITHOUT ATTRIBUTE-TEXT GENERATOR (ATG) ON TRAINING ORDER-1.

| Method | Seen_Avg | | Unseen_Avg | |
|---|---|---|---|---|
| | mAP | R-1 | mAP | R-1 |
| Training w/o ATG | 60.1 | 70.5 | 59.3 | 56.5 |
| Training w/ ATG | **61.8** | **71.9** | **60.8** | **58.3** |

DCR model can effectively bridge domain gaps, improving knowledge transfer capabilities, benefiting from the attribute-oriented anti-forgetting (AF) strategy based on attribute-wise representations.

**The effectiveness of maximizing intra-domain discrimination.** We visualize the feature distribution of KRKC and our method. Figure 4 shows that our DCR model can significantly cluster images of the same identity more tightly (circle) and increase the distance between different identities (black bidirectional arrow). Compared to KRKC, our DCR model effectively improves intra-domain discrimination, benefiting from the complementary relationship between global and attribute representations, which allows it to learn the subtle nuances of individuals.

**Generalization Curves on Unseen Dataset.** We analyze the average performance on unseen datasets over the training steps, as shown in Figure 5. Compared to other methods, our DCR model achieves superior performance and exhibits faster performance growth across the training steps. Thus, our attribute-oriented anti-forgetting (AF) strategy effectively bridges inter-domain gaps, enhancing the generalization ability of our model. In summary, our DCR model explores global and attribute-wise representations to achieve a trade-off between maximizing intra-domain discrimination and minimizing inter-domain gaps.

### C. Ablation Studies

**The number of global and attribute-wise representations.** Global and attribute-wise representations capture individual nuances in intra-domain and inter-domain consistency. We study the suitability of multiple global and attribute-wise representations, as shown in Table III. We observe that setting the number of global and attribute-wise representations $N$ to 3 achieves the best performance for our method.

**Performance of Different Components.** To assess the contribution of each component to our DCR, we conduct ablation studies on seen and unseen datasets, as shown in Table IV. Comparing the first and second rows, we observe that the parallel fusion module (PFM), which employs a parallel cross-attention mechanism, effectively fuses text and image embeddings. Comparing the second and fourth rows, we consider that the attribute compensation network (ACN) and attribute-oriented anti-forgetting (AF) strategy effectively learn domain consistency, improving generalization ability. In the second and third rows, we notice the performance decrease when using only the knowledge consolidation (KC) strategy based on global representations across increasing data while ignoring inter-domain gaps. The results demonstrate that both global representations and attribute-wise representations achieve a trade-off between maximizing intra-domain discrimination and minimizing inter-domain gaps for enhancing the anti-forgetting and generalization capacity of our DCR.

**Performance of attribute-text generator.** To better understand whether each instance's text descriptions generated by the attribute-text generator (ATG) provide more fine-grained guidance for learning global representations, we train our model using the generic text descriptor "A photo of a person" (w/o ATG) for comparison. Table V shows that the attribute-text generator obtain text descriptions to significantly improves overall performance. When using the specific text descriptors, the average decreases by 1.7% mAP/1.4% R-1 on seen datasets and by 1.5% mAP/1.8% R-1 on unseen datasets. ATG enhances the robustness of global representations for each instance, effectively mitigating the forgetting of old knowledge.

## IV. CONCLUSION

In this paper, we propose a domain consistency representation learning (DCR) model that explores global and attribute-wise representations to capture subtle nuances in intra-domain and inter-domain consistency, achieving a trade-off between maximizing intra-domain discrimination and minimizing inter-domain gaps. Specifically, global and attribute-wise representations serve as complementary information to distinguish similar identities in intra-domain. We further develop an attribute-oriented anti-forgetting (AF) strategy and a knowledge consolidation (KC) strategy to minimize inter-domain gaps and facilitate knowledge transfer, enhancing generalization capabilities. Extensive experiments demonstrate that our method achieves superior performance compared to state-of-the-art LReID methods.

## REFERENCES

[1] Wenhang Ge, Junlong Du, Ancong Wu, Yuqiao Xian, Ke Yan, Feiyue Huang, and Wei-Shi Zheng. Lifelong person re-identification by pseudo task knowledge preservation. In *AAAI*, volume 36, pages 688–696, 2022.

[2] Chunlin Yu, Ye Shi, Zimo Liu, Shenghua Gao, and Jingya Wang. Lifelong person re-identification via knowledge refreshing and consolidation. In *AAAI*, volume 37, pages 3295–3303, 2023.

[3] Anguo Zhang, Yueming Gao, Yuzhen Niu, Wenxi Liu, and Yongcheng Zhou. Coarse-to-fine person re-identification with auxiliary-domain classification and second-order information bottleneck. In *CVPR*, pages 598–607, 2021.

[4] Huijie Fan, Xiaotong Wang, Qiang Wang, Shengpeng Fu, and Yandong Tang. Skip connection aggregation transformer for occluded person reidentification. *IEEE Transactions on Industrial Informatics*, 20(1):442–451, 2023.

[5] Siyuan Li, Li Sun, and Qingli Li. Clip-reid: exploiting vision-language model for image re-identification without concrete text labels. In *AAAI*, volume 37, pages 1405–1413, 2023.

[6] Zexian Yang, Dayan Wu, Chenming Wu, Zheng Lin, Jingzi Gu, and Weiping Wang. A pedestrian is worth one prompt: Towards language guidance person re-identification. In *CVPR*, pages 17343–17353, 2024.

[7] Tao Wang, Hong Liu, Pinhao Song, Tianyu Guo, and Wei Shi. Pose-guided feature disentangling for occluded person re-identification based on transformer. In *AAAI*, volume 36, pages 2540–2549, 2022.

[8] Nan Pu, Yu Liu, Wei Chen, Erwin M Bakker, and Michael S Lew. Meta reconciliation normalization for lifelong person re-identification. In *ACM MM*, pages 541–549, 2022.

[9] Jinze Huang, Xiaohan Yu, Dong An, Yaoguang Wei, Xiao Bai, Jin Zheng, Chen Wang, and Jun Zhou. Learning consistent region features for lifelong person re-identification. *Pattern Recognition*, 144:109837, 2023.

[10] Yuming Yan, Huimin Yu, Yubin Wang, Shuyi Song, Weihu Huang, and Juncan Jin. Unified stability and plasticity for lifelong person re-identification in cloth-changing and cloth-consistent scenarios. *IEEE Transactions on Circuits and Systems for Video Technology*, 2024.

[11] Guile Wu and Shaogang Gong. Generalising without forgetting for lifelong person re-identification. In *AAAI*, volume 35, pages 2889–2897, 2021.

[12] Lei Zhang, Guanyu Gao, and Huaizheng Zhang. Spatial-temporal federated learning for lifelong person re-identification on distributed edges. *IEEE Transactions on Circuits and Systems for Video Technology*, 2023.

[13] Kunlun Xu, Xu Zou, and Jiahuan Zhou. Lstkc: Long short-term knowledge consolidation for lifelong person re-identification. In *AAAI*, volume 38, pages 16202–16210, 2024.

[14] Zhicheng Sun and Yadong Mu. Patch-based knowledge distillation for lifelong person re-identification. In *ACM MM*, pages 696–707, 2022.

[15] Kunlun Xu, Xu Zou, Yuxin Peng, and Jiahuan Zhou. Distribution-aware knowledge prototyping for non-exemplar lifelong person re-identification. In *CVPR*, pages 16604–16613, 2024.

[16] Zhenyu Cui, Jiahuan Zhou, Xun Wang, Manyu Zhu, and Yuxin Peng. Learning continual compatible representation for re-indexing free lifelong person re-identification. In *CVPR*, pages 16614–16623, 2024.

[17] Nan Pu, Wei Chen, Yu Liu, Erwin M Bakker, and Michael S Lew. Lifelong person re-identification via adaptive knowledge accumulation. In *CVPR*, pages 7901–7910, 2021.

[18] Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.

[19] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022.

[20] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *CVPR*, pages 16816–16825, 2022.

[21] 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.

[22] Chenliang Li, Haiyang Xu, Junfeng Tian, Wei Wang, Ming Yan, Bin Bi, Jiabo Ye, He Chen, Guohai Xu, Zheng Cao, et al. mplug: Effective and efficient vision-language learning by cross-modal skip-connections. In *Empirical Methods in Natural Language Processing*, pages 7241–7259, 2022.

[23] Chenyang Yu, Xuehu Liu, Yingquan Wang, Pingping Zhang, and Huchuan Lu. Tf-clip: Learning text-free clip for video-based person re-identification. In *AAAI*, volume 38, pages 6764–6772, 2024.

[24] 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 *ICCV*, pages 11174–11184, 2023.

[25] Yang Qin, Yingke Chen, Dezhong Peng, Xi Peng, Joey Tianyi Zhou, and Peng Hu. Noisy-correspondence learning for text-to-image person re-identification. In *CVPR*, pages 27197–27206, 2024.

[26] Xinyi Wu, Wentao Ma, Dan Guo, Tongqing Zhou, Shan Zhao, and Zhiping Cai. Text-based occluded person re-identification via multi-granularity contrastive consistency learning. In *AAAI*, volume 38, pages 6162–6170, 2024.

[27] Yunhao Du, Zhicheng Zhao, and Fei Su. Yyds: Visible-infrared person re-identification with coarse descriptions. *arXiv preprint arXiv:2403.04183*, 2024.

[28] Jian Jia, Houjing Huang, Xiaotang Chen, and Kaiqi Huang. Rethinking of pedestrian attribute recognition: A reliable evaluation under zero-shot pedestrian identity setting. *arXiv preprint arXiv:2107.03576*, 2021.

[29] Haoyun Sun, Hongwei Zhao, Weishan Zhang, Liang Xu, and Hongqing Guan. Adaptive multi-task learning for multi-par in real-world. *IEEE Journal of Radio Frequency Identification*, 2024.

[30] Yunfei Zhou and Xiangrui Zeng. Towards comprehensive understanding of pedestrians for autonomous driving: Efficient multi-task-learning-based pedestrian detection, tracking and attribute recognition. *Robotics and Autonomous Systems*, 171:104580, 2024.

[31] Xinwen Fan, Yukang Zhang, Yang Lu, and Hanzi Wang. Parformer: Transformer-based multi-task network for pedestrian attribute recognition. *IEEE Transactions on Circuits and Systems for Video Technology*, 34(1):411–423, 2023.

[32] Jian Jia, Naiyu Gao, Fei He, Xiaotang Chen, and Kaiqi Huang. Learning disentangled attribute representations for robust pedestrian attribute recognition. In *AAAI*, volume 36, pages 1069–1077, 2022.

[33] Jian Jia, Xiaotang Chen, and Kaiqi Huang. Spatial and semantic consistency regularizations for pedestrian attribute recognition. In *ICCV*, pages 962–971, 2021.

[34] Xiaoyan Yu, Neng Dong, Liehuang Zhu, Hao Peng, and Dapeng Tao. Clip-driven semantic discovery network for visible-infrared person re-identification. *arXiv preprint arXiv:2401.05806*, 2024.

[35] Xihui Liu, Haiyu Zhao, Maoqing Tian, Lu Sheng, Jing Shao, Shuai Yi, Junjie Yan, and Xiaogang Wang. Hydraplus-net: Attentive deep features for pedestrian analysis. In *ICCV*, pages 350–359, 2017.

[36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[37] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.

[38] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *ICCV*, pages 1116–1124, 2015.

[39] Tong Xiao, Shuang Li, Bochao Wang, Liang Lin, and Xiaogang Wang. End-to-end deep learning for person search. *arXiv preprint arXiv:1604.01850*, 2(2):4, 2016.

[40] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi. Performance measures and a data set for multi-target, multi-camera tracking. In *ECCV*, pages 17–35. Springer, 2016.

[41] Longhui Wei, Shiliang Zhang, Wen Gao, and Qi Tian. Person transfer gan to bridge domain gap for person re-identification. In *CVPR*, pages 79–88, 2018.

[42] Wei Li, Rui Zhao, Tong Xiao, and Xiaogang Wang. Deepreid: Deep filter pairing neural network for person re-identification. In *CVPR*, pages 152–159, 2014.

[43] Douglas Gray and Hai Tao. Viewpoint invariant pedestrian recognition with an ensemble of localized features. In *ECCV*, pages 262–275. Springer, 2008.

[44] Chen Change Loy, Tao Xiang, and Shaogang Gong. Time-delayed correlation analysis for multi-camera activity understanding. *International Journal of Computer Vision*, 90:106–129, 2010.

[45] Wei Li and Xiaogang Wang. Locally aligned feature transforms across views. In *CVPR*, pages 3594–3601, 2013.

[46] Jiaxu Miao, Yu Wu, Ping Liu, Yuhang Ding, and Yi Yang. Pose-guided feature alignment for occluded person re-identification. In *ICCV*, pages 542–551, 2019.

[47] Jiaxuan Zhuo, Zeyu Chen, Jianhuang Lai, and Guangcong Wang. Occluded person re-identification. In *ICME*, pages 1–6. IEEE, 2018.

[48] Martin Hirzer, Csaba Beleznai, Peter M Roth, and Horst Bischof. Person re-identification by descriptive and discriminative classification. In *Image Analysis: 17th Scandinavian Conference*, pages 91–102. Springer, 2011.

[49] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[50] James Seale Smith, Leonid Karlinsky, Vyshnavi Gutta, Paola Cascante-Bonilla, Donghyun Kim, Assaf Arbelle, Rameswar Panda, Rogerio Feris, and Zsolt Kira. Codaprompt: Continual decomposed attention-based prompting for rehearsal-free continual learning. In *CVPR*, pages 11909–11919, 2023.