

Exploring Shape Embedding for Cloth-Changing Person Re-Identification via 2D-3D Correspondences

Yubin Wang

College of Information Science and
Electronic Engineering, Zhejiang
University
Hangzhou, China
zjuwyb1999@gmail.com

Huimin Yu*

College of Information Science and
Electronic Engineering, Zhejiang
University
Hangzhou, China
yhm2005@zju.edu.cn

Yuming Yan

College of Information Science and
Electronic Engineering, Zhejiang
University
Hangzhou, China
12231016@zju.edu.cn

Shuyi Song

College of Information Science and
Electronic Engineering, Zhejiang
University
Hangzhou, China
22131091@zju.edu.cn

Biyang Liu

College of Information Science and
Electronic Engineering, Zhejiang
University
Hangzhou, China
11831033@zju.edu.cn

Yichong Lu

College of Information Science and
Electronic Engineering, Zhejiang
University
Hangzhou, China
luyi200106@gmail.com

ABSTRACT

Cloth-Changing Person Re-Identification (CC-ReID) is a common and realistic problem since fashion constantly changes over time and people's aesthetic preferences are not set in stone. While most existing cloth-changing ReID methods focus on learning cloth-agnostic identity representations from coarse semantic cues (e.g. silhouettes and part segmentation maps), they neglect the continuous shape distributions at the pixel level. In this paper, we propose Continuous Surface Correspondence Learning (CSCL), a new shape embedding paradigm for cloth-changing ReID. CSCL establishes continuous correspondences between a 2D image plane and a canonical 3D body surface via pixel-to-vertex classification, which naturally aligns a person image to the surface of a 3D human model and simultaneously obtains pixel-wise surface embeddings. We further extract fine-grained shape features from the learned surface embeddings and then integrate them with global RGB features via a carefully designed cross-modality fusion module. The shape embedding paradigm based on 2D-3D correspondences remarkably enhances the model's global understanding of human body shape. To promote the study of ReID under clothing change, we construct 3D Dense Persons (DP3D), which is the first large-scale cloth-changing ReID dataset that provides densely annotated 2D-3D correspondences and a precise 3D mesh for each person image, while containing diverse cloth-changing cases over all four seasons. Experiments on both cloth-changing and cloth-consistent ReID benchmarks validate the effectiveness of our method. Our project page is located at <https://CSCL-CC.github.io>.

*Corresponding Author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).

MM '23, October 29–November 3, 2023, Ottawa, ON, Canada

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0108-5/23/10...\$15.00

<https://doi.org/10.1145/3581783.3611715>

CCS CONCEPTS

• **Computing methodologies** → **Object identification; Object recognition.**

KEYWORDS

Cloth-Changing Person Re-Identification; Shape Embedding; 2D-3D Correspondences; Large-Scale Dataset; Cross-Modality Fusion.

ACM Reference Format:

Yubin Wang, Huimin Yu, Yuming Yan, Shuyi Song, Biyang Liu, and Yichong Lu. 2023. Exploring Shape Embedding for Cloth-Changing Person Re-Identification via 2D-3D Correspondences. In *Proceedings of the 31st ACM International Conference on Multimedia (MM '23)*, October 29–November 3, 2023, Ottawa, ON, Canada. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3581783.3611715>

1 INTRODUCTION

Person Re-Identification (Re-ID) targets at re-identifying a specific person across disjoint cameras [28]. Most existing works [6, 17, 20, 23, 46, 51, 56] presuppose that the appearances of people remain consistent over time. In reality, people tend to change their outfits over a long duration and different people may share the same dressing sense. Methods that rely excessively on clothing appearance fail to generalize to this long-term cloth-changing scenario.

In recent years, plenty of efforts [3, 19, 21, 36, 47, 50, 53] have been made to handle the cloth-changing issue by learning discriminative cloth-agnostic identity representations. A small proportion of methods [3, 47, 50] attempt to decouple cloth-agnostic features directly from RGB images without multi-modal auxiliary information, which inevitably leads to the loss of crucial information in global features and results in a heavy reliance on the domain. The mainstream methods [7, 19, 21, 29, 36, 53] typically adopt human parsing models to obtain coarse semantic cues to guide the extraction of biometric features, such as shape features. However, as shown in Figure 1(b), coarse semantic cues are insufficient to obtain detailed shape information of a specific person, as it only enables the estimation of body part labels but fails to model pixel-wise shape distributions within the parts. Several recent works [22, 54]

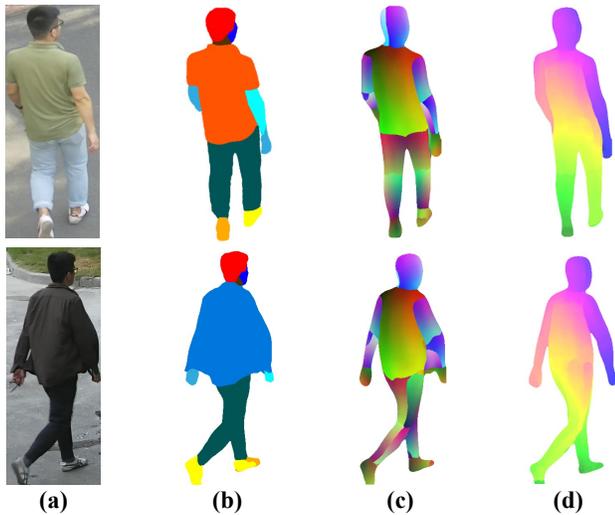


Figure 1: Comparison of different multi-modal auxiliary information for person re-identification. (a) Images of the same person in DP3D; (b) Coarse part segmentation, with only part labels estimated; (c) Discretized DensePose IUUV estimation, with obvious seams between body parts; (d) Continuous 2D-3D correspondences between image pixels and the entire body surface, obtained through our CSCL framework.

leverage dense pose estimation [1] to align the texture of body parts based on UV mapping. However, they do not further explore reliable shape representations for the ReID task. Additionally, these methods have a major defect in that they require partitioning the 3D model into charts, and the resulting discretized UV spaces prevent them from learning continuous correspondences over the entire body surface. As shown in Figure 1(c), the use of independent UV coordinate systems for each body part results in noticeable part seams in the estimated IUUV maps. There are also some methods [4, 43] directly estimating SMPL [27] shape parameters as 3D shape features. However, the SMPL shape parameter space is highly incompatible with the image feature space, making it challenging to effectively integrate features from these two modalities.

In this paper, we propose a Continuous Surface Correspondence Learning (CSCL) framework, which represents a new shape embedding paradigm for cloth-changing ReID. CSCL pixel-wisely maps a person image to a continuous embedding space of the SMPL mesh surface through vertex classification. Essentially, learning continuous 2D-3D correspondences aligns a person image to the entire surface of a 3D human model, and simultaneously obtains a pixel-level continuous distribution of body shape on the canonical 3D surface. Even for different persons wearing the same clothes, there can be significant differences in their body shape distributions. Therefore, we further extract fine-grained discriminative shape features from the established correspondences, and integrate them with global RGB features via an optimized cross-modality fusion module based on the transformer [39], which greatly compensates for the lost shape details in global RGB features. We incorporate a novel Latent Convolutional Projection (LCP) layer for feature projection. The LCP layer enhances the sharing and correlation

among tokens via adding an additional latent embedding, which is the latent vector of an auto-encoder designed to reconstruct the token map. It is also noteworthy that the proposed framework generalizes well to the cloth-consistent cases, indicating the reliability of the learned shape features.

However, there is currently no publicly available cloth-changing ReID dataset with ground-truth dense 2D-3D correspondences. To facilitate the research, we construct a large-scale cloth-changing ReID dataset named 3D Dense Persons (DP3D), which contains 39,100 person images of 413 different persons captured by 15 cameras over all four seasons. We annotated dense 2D-3D correspondences for each person image via a carefully designed annotation system, ensuring 80 to 125 annotations for each image.

The main contributions of this work are summarized as follows:

- We propose a new shape embedding paradigm for cloth-changing ReID that establishes pixel-wise and continuous correspondences between a 2D image plane and a canonical 3D human body surface. To the best of our knowledge, this is also the first work to explore global shape representations for cloth-changing ReID via 2D-3D correspondences.
- We develop an optimized cross-modality fusion module to adaptively integrate shape features with global RGB features, where a novel Latent Convolutional Projection (LCP) layer is designed to perform feature projection.
- We construct 3D Dense Persons (DP3D), which is the first large-scale cloth-changing ReID dataset with densely annotated 2D-3D correspondences and a corresponding 3D mesh for each person image, while containing highly diverse cloth-changing cases in real-world scenarios.
- We demonstrate our proposed method is applicable to both cloth-changing and cloth-consistent situations, as shown by extensive results on four cloth-changing ReID datasets including DP3D and two general ReID datasets.

2 RELATED WORKS

In this section, we first review the literature on cloth-changing person re-identification and corresponding datasets, then introducing the research related to continuous surface embeddings in the context of 3D shape analysis.

2.1 Cloth-Changing Person ReID

Existing cloth-changing ReID methods can be categorized into decoupling-based methods and auxiliary modality-based methods. Decoupling-based methods [9, 47, 48] aim to decouple cloth-agnostic features directly from RGB images without multi-modal auxiliary information. AFD-Net [47] disentangled identity and clothing features via generative adversarial learning. CAL [9] proposed to penalize the predictive power of the ReID model with respect to clothes via a clothes-based adversarial loss, while UCAD [48] enforced the identity and clothing features to be linearly independent in the feature space via an orthogonal loss.

Auxiliary modality-based methods [4, 7, 15, 22, 36, 54] are considered more robust since visual texture features can be filtered under the supervision of human semantics. FSAM [15] proposed to complement 2D shape representations obtained from human silhouettes for global features. MVSE [7] embedded multigranular

Table 1: Comparison of DP3D and existing cloth-changing ReID datasets ('In': Indoor; 'Out': Outdoor).

Datasets	Scene	IDs	Image	Cam	Time	3D View	Dense Corr.
Celebrities [16]	-	590	10,842	-	-	✗	✗
LTCC [33]	In	152	17,138	12	2 Months	✗	✗
PRCC [49]	In	221	33,698	3	-	✗	✗
COCAS [52]	In	5,266	62,382	30	-	✗	✗
VC-Clothes [40]	-	512	19,060	-	-	✗	✗
CSCC [48]	Out	267	36,700	13	12 Months	✗	✗
NKUP [42]	In/Out	107	9,738	15	4 Month	✗	✗
NKUP+ [26]	In/Out	361	40,217	29	10 Month	✗	✗
DP3D (Ours)	Out	413	39,100	15	12 Months	✓	✓

visual semantic information into the model. Pixel Sampling [36] leveraged a human parsing model to recognize upper clothes and pants, and then randomly changed them by sampling pixels from other people, enforcing the model to automatically learn cloth-agnostic cues. DSA-ReID[54] and ASAG-Net[22] proposed to use dense human semantics to generate semantics-aligned images in the discretized DensePose UV space, while 3DSL [4] considered the low-dimensional SMPL shape parameters as 3D shape features, and directly fused them to global features. None of these methods consider establishing pixel-wise and continuous 2D-3D correspondences between image pixels and the entire 3D body surface, which effectively bridges the gap between 2D and 3D shape space.

2.2 Cloth-Changing ReID Datasets

General person ReID datasets[24, 34, 44, 55] assume that the appearance of the same individual is consistent, which is often not the case in real-world scenarios. Models trained on these datasets rely excessively on clothing appearance, making it difficult for them to generalize well to long-term cloth-changing scenarios. In recent years, a few datasets were collected specifically for the cloth-changing setting. Celebrities [16] were obtained from the Internet, which consists of street snapshots of celebrities. PRCC [49] provides indoor cloth-changing person images with their corresponding contour sketches. COCAS [52] is a large-scale dataset that provides a variety of clothes templates for cloth-changing person ReID. LTCC [33] assumes that different people wear different clothes and assigns a unique clothing label to each person image in the dataset. VC-Clothes [40] is a large realistic synthetic dataset rendered by the GTA5 game engine. CSCC [48] considers different degrees of cloth-changing. NKUP [42] contains both indoor and outdoor person images with complex illumination conditions, while NKUP+ [26] has more diverse scenarios, perspectives, and appearances.

2.3 Continuous Surface Embeddings

Continuous Surface Embeddings (CSE) target at pixel-wisely learning an embedding of the corresponding 3D vertex from an RGB image [30], which demonstrates strong human body representation capabilities. HumanGPS [38] employs contrastive learning to enhance CSE representations. BodyMap [18] introduced a coarse-to-fine learning scheme, establishing high-definition full-body continuous correspondences by refining coarse correspondences. SurfEmb [10] applied Continuous Surface Embeddings to the field of object pose estimation and learned correspondence distributions in a self-supervised fashion.

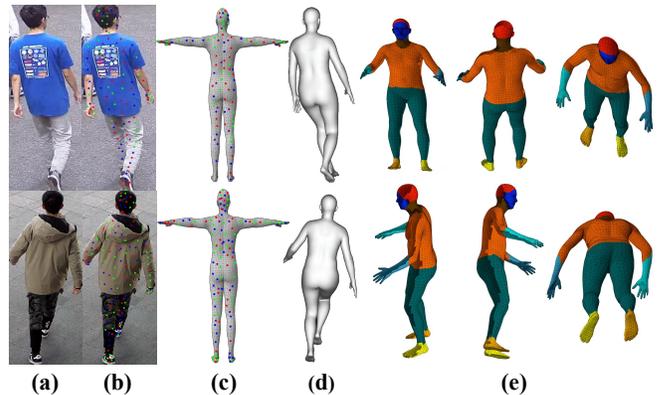


Figure 2: Examples of annotating person images in the DP3D dataset. (a) Cross-appearance images of the same person; (b) Generating pixels to be labeled (corresponding pixels are visualized with purple dots); (c) Annotating ground-truth corresponding 3D mesh vertices. (d) fitting the SMPL model to the person images under the guidance of dense correspondences; (e) the projected 2D full-body images used for annotation.

3 THE 3D DENSE PERSONS DATASET

Obtaining ground-truth 3D structure information for pedestrians is of substantial importance as it can address potential geometric ambiguities that may arise from relying solely on RGB modality.

In this section, we introduce the 3D Dense Persons (DP3D), a large-scale cloth-changing ReID dataset that provides densely annotated 2D-3D correspondences and a corresponding 3D mesh for each person image, filling the gap in the field.

3.1 Data Collection

The raw videos we collected have high resolutions and cover a time span of one year. We selected a total of 15 cameras, with 5 of them having a resolution of 4K, 2 having a resolution of 2K, and the remainder being set to a resolution of 1080P. The use of high-resolution cameras ensures the recorded pedestrians to be as clear as possible, which is advantageous for the ReID task under clothing change. The shooting scenes encompass various outdoor locations, such as street scenes, park landscapes, construction sites, and parking lots. All pedestrians were captured by at least 2 cameras, with the majority being captured by 3 or more. We adopted the Mask R-CNN [11] framework to detect the bounding box of each person after framing.

3.2 Annotation System

Due to the dramatic variations in people’s clothing styles over the course of a year, we first identified the volunteers and conducted a manual inspection to avoid misidentification, while assigning a camera ID label, a person ID label, and a clothing ID label to each person image. Then, as shown in Figure 2, we annotated dense correspondences via a carefully designed pipeline. In the first stage, we ran the universal model of Graphonomy [8] with 20 part labels to segment the images, then uniformly sampling 40 pixels across the entire human body region. We also utilized

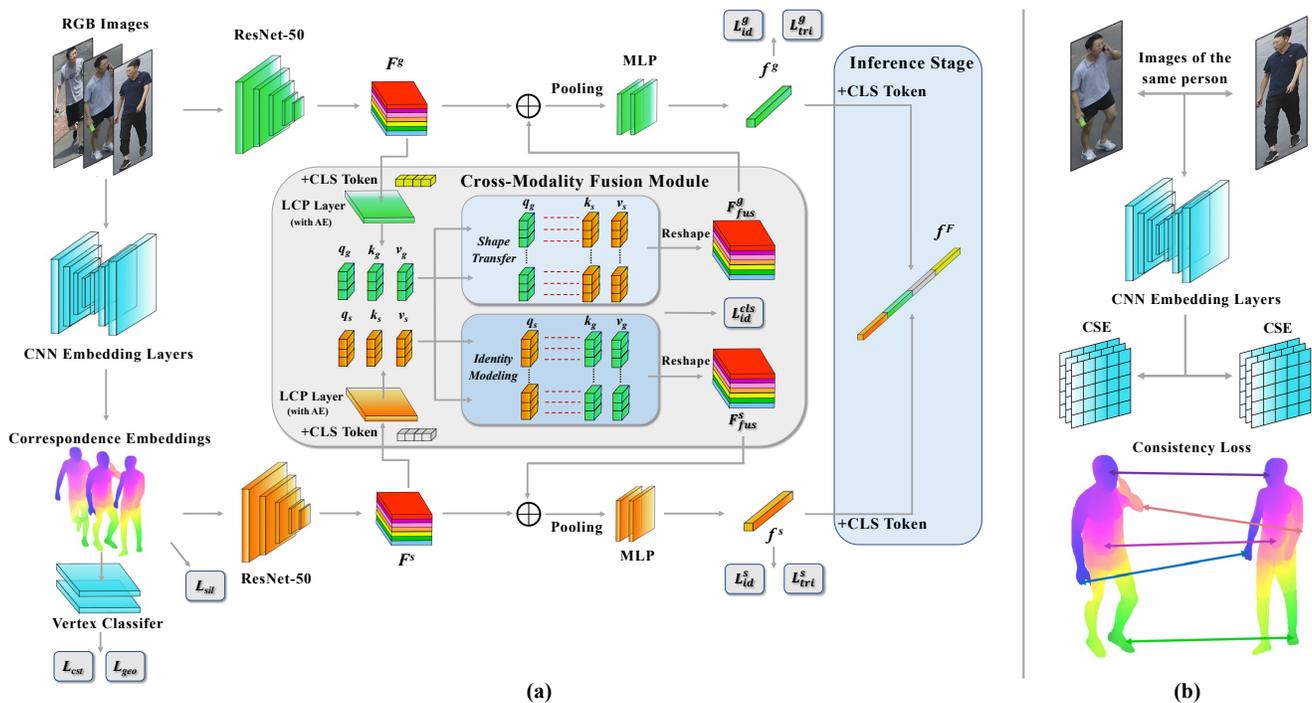


Figure 3: The architecture of the CSCL framework. (a) Our framework learns pixel-wise and continuous 2D-3D correspondences, which enables the extraction of fine-grained shape features. Cloth-agnostic shape knowledge is then complemented for global RGB features via cross-modality fusion; (b) Consistency learning between cross-view corresponding pixels.

k-means clustering to obtain 5 to 10 centroid pixels for each part based on its size. Compared to DensePose [1], our sampling method avoids seams between body parts and ensures a sufficient number of sampling points for smaller parts. However, since people may wear loose clothes, we manually filtered out those sampling pixels that did not fall within the human body regions underneath the clothes. For each pair of images belonging to the same person, we additionally selected 10 corresponding pixels for consistency learning, which correspond to the same 10 mesh vertices. In the second stage, as shown in Figure 2 (e), we projected the SMPL mean template mesh from 6 predefined viewpoints to generate full-body images. When annotating a specific pixel, it was only necessary to choose the most suitable projected image, and its 2D coordinates were used to localize the corresponding 3D vertex. In cases certain pixels were challenging to determine from the projected images, we directly annotated the correspondences on the 3D mesh surface through rotation. It is worth noting that we did not annotate in a part-by-part manner, but rather adopted a global approach using full-body projected images for annotation, which ensured accurate annotations at the junctions of body parts. In the last stage, to obtain accurate SMPL parameters, we employed a modified SMPLity-X [32] to fit the SMPL model to the person images under the guidance of densely annotated correspondences.

3.3 Statistics and Comparison

The proposed DP3D dataset is characterized by its diverse scenes, multiple perspectives, large number of individuals, and long time

span. It comprises 39,100 person images belonging to 413 different persons, which were captured over the course of a year (during four distinct seasons). Depending on its resolution, each person image has approximately 80 to 125 annotated correspondences, where 10 correspondences have mesh vertices shared among all images of the same person. We divided the images into a training set and a testing set, with each set containing approximately equal numbers of identities. For same-appearance images of a specific person, we randomly select one image per viewpoint to construct the query set, while the remaining images in the testing set form the gallery set. We present in Table 1 a comparison between DP3D and existing cloth-changing ReID datasets.

4 METHODOLOGY

In this section, we first provide an overview of our proposed framework in Section 4.1. Next, in Section 4.2 and 4.3, we elaborate the learning scheme of continuous 2D-3D correspondences, as well as the design principles of the cross-modality fusion module, respectively. Subsequently, we provide a comprehensive description of the training losses in Section 4.4.

4.1 Overview

As shown in Figure 3 (a), person images are input separately into the ResNet-50 [12] backbone and CNN embedding layers to extract global RGB features and continuous surface embeddings. For each foreground pixel, CSCL maps it to a continuous embedding space of the SMPL mesh surface under the supervision of geodesic distances.

Subsequently, a shape extraction network with a ResNet-50 architecture is further employed to extract fine-grained shape features from the learned surface embeddings, while simultaneously mapping them to the same size as global RGB features. Following that, we adaptively integrate shape features with global RGB features via an improved cross-modality fusion module, where a novel Latent Convolutional Projection (LCP) layer is designed to perform feature projection. Cross-attention mechanism is then applied to aggregate features from the two distinct modalities, which are then added to the original features. After the fusion, we conduct Global Average Pooling (GAP), followed by two separate fully-connected classifiers, to obtain the final global RGB features and shape features. We also introduce a learnable class token for each of the two modalities, which exhibits strong cross-modality compatibility and also contributes to the ID loss. In the inference stage, the two class tokens are concatenated with global RGB features and shape features to construct the final identity feature.

4.2 Establishing Continuous Correspondences

Considering the huge domain gap between 2D person images and the 3D space perceived by human eyes, we believe that establishing continuous correspondences between image pixels and the entire 3D human body is of substantial importance, which bridges the gap between the 2D and 3D shape space and therefore benefit the understanding of global body shape.

Given a person image $I \in \mathbb{R}^{H \times W \times 3}$ of height H and width W , we first extract the segmentation mask M of the foreground person. Then, the CNN embedding layers map the person image into continuous surface embeddings $E \in \mathbb{R}^{H \times W \times D}$, while preserving the spatial resolution of the image. For pixels within the foreground mask M , we employ geodesic distances on the 3D surface to supervise the learning of surface embeddings. More concretely, we scale the cross-entropy loss of pixel-to-vertex classification on the mesh surface using geodesic distances. This constraint is reasonable as it quantifies the deviation of vertex prediction on the 3D surface. Furthermore, as illustrated in Figure 3 (b), we also conduct consistency learning for corresponding pixels in images that belong to the same person. Suppose we have two distinct images of the same person, denoted as I_1, I_2 , where foreground pixels p_1 and p_2 belong to image I_1 , and pixel q belongs to image I_2 . Both p_1 and q correspond to the same vertex v_1 on the mesh surface, while p_2 corresponds to vertex v_2 . We first compute the cosine distance in the embedding space to measure the similarity between p_1 and q :

$$d(p_1, q) = 1 - \cos(E_1(p_1), E_2(q)) \quad (1)$$

where E_1 and E_2 denote surface embeddings of images I_1 and I_2 . By minimizing the cosine distance $d(p_1, q)$, the embedding vectors of two corresponding pixels are brought closer. However, during training, only considering the consistency of corresponding pixels may lead to all embeddings mapping to similar values. Therefore, for different pixels p_1 and p_2 in the same person image, we keep their relative affinity by enforcing embedding distances to follow geodesic distances, i.e. minimizing $|d(p_1, p_2) - s(g(v_1, v_2))|$, where $g(\cdot, \cdot)$ calculates the geodesic distance between two mesh vertices and $s(\cdot)$ scales it to match the range of the cosine distance $d(\cdot, \cdot)$.

Establishing 2D-3D correspondences allows for learning the continuous shape distributions on the 3D surface at the pixel level,

i.e. $Pr(v|I, p, p \in M)$, where I denotes the person image, and M denotes the foreground mask. To further extract fine-grained shape features, we feed the learned embeddings into the shape extraction network with a ResNet-50 architecture, while mapping them to the same size as global RGB features. Note that the extracted shape features are insensitive to clothing appearance as texture features are already filtered out in the correspondence learning process.

4.3 Cross-Modality Feature Fusion

To adaptively integrate the shape features extracted from the established continuous correspondences with global RGB features, a cross-modality fusion module is designed. As discussed in CVT [45], convolutional layers are renowned for their remarkable ability to capture intricate local spatial token structures, which allows the removal of positional embeddings from the transformer [39] framework. However, the utilization of fixed-size convolutional kernels hampers the effectiveness of capturing global positional correlations between non-adjacent tokens. To mitigate this issue, we propose a novel Latent Convolutional Projection (LCP) layer. It adds the same latent embedding to each token in the token map, which is the latent vector of a pretrained auto-encoder designed to reconstruct the token map. During the training of CSCL, only the encoder of the auto-encoder is preserved and fixed to ensure the universal nature of the latent embedding, whereas the decoder is disregarded. This design not only greatly enhances the correlation and sharing among different tokens, but also enables better adaptation to images with diverse backgrounds. The projection of an LCP layer can be formulated as follows:

$$Q/K/V = \text{Flatten}(\text{Conv2d}(\text{Reshape2D}(F) + l)) \quad (2)$$

where $Q/K/V$ represents the projected queries, keys, and values, F is the input token map, l represents the latent embedding, and Reshape2D denotes the operation to reshape the feature map F to a 2D token map. After separately passing global RGB features $F^g \in \mathbb{R}^{h \times w \times c}$ and shape features $F^s \in \mathbb{R}^{h \times w \times c}$ through two distinct LCP layers, the cross-attention mechanism is applied to adaptively integrate features from different modalities. We first take global RGB features as queries and shape features as keys/values, reshape the fused feature to match the size of F^g , and finally add it to F^g :

$$F^g = F^g + \text{Reshape3D}(\text{MHA}(Q_g, K_s, V_s)) \quad (3)$$

where Reshape3D denotes the operation of reshaping a 2D token map to match the size of F^g , and MHA represents the multi-head attention. We also take shape features as queries and global RGB features as keys/values for identity modeling of shape features.

$$F^s = F^s + \text{Reshape3D}(\text{MHA}(Q_s, K_g, V_g)) \quad (4)$$

In other words, we enable bidirectional access between global RGB features and shape features, which allows the model not only complements fine-grained cloth-agnostic shape knowledge for global RGB features F^g , but also integrates essential identity-related characteristics for shape features F^s to assist identity modeling. Additionally, we introduce learnable class tokens for each of the two modalities, which are also utilized to compute the ID loss.

4.4 Loss Function

CSE Losses. As discussed in Section 4.2, to mask out the background pixels, the foreground silhouette for each person image is retrieved, thus a binary cross-entropy loss \mathcal{L}_{sil} is employed to penalize unsatisfactory silhouette predictions. Furthermore, we employ geodesic distances on the mesh surface to scale the per-pixel vertex classification loss, which penalizes the misclassified pixels based on the degree of deviation on the surface. The geodesic loss can be formulated as follows:

$$\mathcal{L}_{geo} = -\frac{1}{N} \sum_{p \in I} g(v_p, \hat{v}_p) \cdot \log(p(\hat{v}_p)) \quad (5)$$

where N indicates the number of pixels with ground-truth annotations in image I , v_p and \hat{v}_p represent the ground-truth and predicted mesh vertices corresponding to pixel p , and $g(\cdot, \cdot)$ calculates geodesic distances between two mesh vertices. For consistency learning of continuous surface embeddings, we design the following consistency loss \mathcal{L}_{cst} :

$$\begin{aligned} \mathcal{L}_{cst} = & \frac{1}{N_1} \sum_{p \in I_1, q \in I_2} \log(1 + \exp(d(p, q))) \\ & + \frac{1}{N_2} \sum_{p_1, p_2 \in I} \log(1 + \exp(|d(p_1, p_2) - s(g(v_1, v_2))|)) \end{aligned} \quad (6)$$

where N_1 and N_2 indicate the number of annotated pairs, p and q are corresponding pixels in cross-view images, p_1 and p_2 stand for different pixels in the same image, $d(\cdot, \cdot)$ and $g(\cdot, \cdot)$ respectively denote the cosine distance in the embedding space and the geodesic distance on the surface, and $s(\cdot)$ represents the scale function. The first term of \mathcal{L}_{cst} ensures consistency between embeddings of cross-view corresponding pixels, while the second term enforces embedding distances to follow geodesic distances for different pixels in the same image, thus pushing apart their embeddings, and avoiding the degradation cases that may occur during training.

ReID Losses. The ReID losses employed in our framework consist of a cross-entropy loss (ID loss) for classification and a triplet loss [14] for similarity learning in the feature space. The final global RGB feature f^g , shape feature f^s , and two class tokens all contribute to the ID loss:

$$\mathcal{L}_{id} = \mathcal{L}_{id}^g + \mathcal{L}_{id}^s + \mathcal{L}_{id}^{cls} \quad (7)$$

where \mathcal{L}_{id}^{cls} represents the summation of ID losses of the two class tokens. We introduce separate triplet losses for global RGB features and shape features to enhance their discriminative capability, which are combined to obtain the final triplet loss:

$$\mathcal{L}_{tri} = \mathcal{L}_{tri}^g + \mathcal{L}_{tri}^s \quad (8)$$

Final Loss. The overall objective function of our proposed Continuous Surface Correspondence Learning (CSCL) framework comprises the aforementioned CSE losses and ReID losses, which can be formulated as follows:

$$\mathcal{L} = \mathcal{L}_{sil} + \lambda_1(\mathcal{L}_{geo} + \alpha \mathcal{L}_{cst}) + \lambda_2 \mathcal{L}_{id} + \lambda_3 \mathcal{L}_{tri} \quad (9)$$

where λ_1 , α , λ_2 and λ_3 are weights for balancing each term.

5 EXPERIMENTS

5.1 Datasets and Protocols

We conduct experiments on four existing cloth-changing ReID datasets (i.e. LTCC [33], PRCC [49], VC-Clothes [40] and DP3D). Furthermore, three different settings are involved in our experiment: (1) **Standard Setting**: the test set includes both same-appearance and cross-appearance samples; (2) **Cloth-Changing Setting**: the test set only includes cross-appearance samples; (3) **Same-Clothes Setting**: the test set only includes same-appearance samples. For LTCC and DP3D, we provide experimental results in the standard setting and cloth-changing setting, while for PRCC and VC-Clothes, results in the same-clothes setting and cloth-changing setting are reported. We additionally validate our method on two general ReID datasets (i.e. Market-1501 [55] and DukeMTMC [34]), following their evaluation metrics. For evaluation, we adopt the mean average precision (mAP) and rank-1 accuracy to evaluate the effectiveness of ReID methods. We also utilize Geodesic Point Similarity (GPS) [1] scores to measure the quality of the established correspondences:

$$GPS_I = \frac{1}{N} \sum_{p \in I} \exp \frac{-g(v_p, \hat{v}_p)^2}{2\sigma^2} \quad (10)$$

where I indicates a person image, N is the number of ground-truth correspondences, v_p and \hat{v}_p denote the ground-truth vertex and the estimated vertex, $g(\cdot, \cdot)$ represents geodesic distances, and σ is a normalizing factor set to 0.255. When GPS scores exceed a certain threshold, the correspondences are considered as correct. Therefore, following the metric of BodyMap [18], we report Average Precision (AP) and Average Recall (AR) based on GPS scores.

5.2 Implementation Details

For datasets without ground-truth dense correspondences, we fit the SMPL body model to the person images under the guidance of OpenPose [2] keypoint detections and foreground silhouettes. For each SMPL mesh vertex, there is a reprojected point on the 2D image plane, and the pixel closest to this point is utilized to establish the correspondence. If different vertices correspond to the same pixel, only the vertex closest to the camera is recorded. Based on the image resolution, we uniformly sampled 80 to 125 pseudo correspondences within the entire body region. All input images are resized to 256×128 . A skip-connecting UNet [35] architecture pretrained on the DensePose-COCO dataset [1] is employed as embedding layers, while two distinct ResNet-50 backbone pretrained on ImageNet [5] with the last downsampling layer discarded are employed to extract global RGB features and shape features, respectively. In the training stage, the Adam optimizer[31] was utilized for optimization. We first trained the embedding layers for 50 epochs with a learning rate of 5×10^{-5} , and then fixed them to train the rest of the network for 100 epochs with a linear warm-up phase. The learning rate was increased from 1×10^{-5} to 1×10^{-4} in the first 5 epochs. Finally, we trained the network in an end-to-end manner for 40 epochs with a fixed learning rate of 1×10^{-5} . The embedding dimension D is set to 64. The values of λ_1 , α , λ_2 , λ_3 in Eq. 9 are set to 0.3, 5.0, 1.0, 0.8, and the margin parameter for the triplet loss is set to 0.3, respectively.

Table 2: Comparison on LTCC, PRCC, VC-Clothes and DP3D datasets. # denotes we conducted experiments based on the code we reproduced. ‘Standard’, ‘Cloth-Changing’ and ‘Same-Clothes’ represent experiment settings illustrated in Section 5.1.

Methods	LTCC				PRCC				VC-Clothes				DP3D			
	Standard		Cloth-Changing		Same-Clothes		Cloth-Changing		Same-Clothes		Cloth-Changing		Standard		Cloth-Changing	
	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP
PCB (ECCV18) [37]	65.1	30.6	23.5	10.0	86.9	83.6	22.9	24.7	72.3	73.9	53.9	55.6	58.3	35.9	15.1	9.9
HACNN (CVPR18) [25]	60.2	26.7	21.5	9.2	82.4	84.7	21.8	23.2	68.6	69.7	49.6	50.1	53.4	31.8	13.4	8.5
MGN (MM18) [41]	68.4	32.4	25.3	11.5	89.8	87.4	25.9	35.9	74.3	75.2	55.0	57.3	59.7	37.0	17.9	12.2
TransReID (ICCV21) [13]	70.1	33.8	26.4	12.6	93.1	94.0	40.1	43.6	79.8	80.3	73.1	74.9	62.5	37.5	18.5	12.7
SE+CESD (ACCV20) [33]	71.4	34.3	26.2	12.4	91.8	90.6	37.6	38.7	85.2	79.1	69.5	65.5	61.9	38.3	18.3	12.7
FSAM (CVPR21) [15]	73.2	35.4	38.5	16.2	98.8	-	54.5	-	94.7	94.8	78.6	78.9	61.7	39.0	17.7	11.9
3DSL (CVPR21) [4]	73.8#	34.2#	31.2	14.8	98.7#	95.0#	51.3	49.8#	92.5#	79.7#	79.9	81.2	66.4#	45.3#	29.6#	17.8#
UCAD (IJCAI22) [48]	74.4	34.8	32.5	15.1	96.5	95.9	45.3	45.2	92.6	81.1	82.4	73.8	63.5	41.7	21.3	13.1
MVSE (MM22) [7]	73.4#	33.9#	70.5	33.0	98.7#	98.3#	47.4	52.5	86.1#	79.5#	79.4#	79.1#	63.7#	41.7#	21.2#	13.4#
M2NET (MM22) [26]	-	-	-	-	99.5	99.1	59.3	57.7	-	-	-	-	63.3	39.4	20.8	12.9
CAL (CVPR22) [9]	74.2	40.8	40.1	18.0	100	99.8	55.2	55.8	-	-	-	-	64.8	42.4	22.9	14.4
Baseline(ResNet-50)	68.2	34.3	26.2	12.3	89.6	88.0	32.8	37.1	78.0	78.8	70.6	65.9	62.5	38.8	19.2	13.0
CSCL(w/o. \mathcal{L}_{cst})	75.3	41.1	68.9	33.5	99.7	99.4	63.5	63.6	97.1	95.4	85.5	84.7	74.1	55.8	37.8	27.0
CSCL	75.5	41.6	69.7	34.1	99.7	99.6	64.2	64.5	97.3	95.5	85.9	84.7	75.8	56.9	39.2	28.7

Table 3: Comparison of CSCL and other competitors on Market-1501 (single-query setting) and DukeMTMC.

Methods	Market-1501		DukeMTMC	
	Rank-1	mAP	Rank-1	mAP
PCB (ECCV18) [37]	92.3	77.4	81.8	66.1
HACNN (CVPR18) [25]	91.2	75.7	80.5	63.8
MGN (MM18) [41]	95.7	86.9	88.7	78.4
Trans-ReID (ICCV21) [13]	95.2	89.5	90.7	82.6
3DSL (CVPR21) [4]	95.0	87.3	88.2	76.1
Baseline(ResNet-50)	92.7	78.0	85.8	75.3
CSCL	95.4	89.5	90.3	83.1

5.3 Comparison with State-of-the-arts

As shown in Table2, we compare our proposed CSCL with seven SOTA cloth-changing methods (i.e. SE+CESD [33], PSAM [15], 3DSL [4], UCAD [48], MVSE [7], M2NET [26] and CAL [9]) on LTCC, PRCC, VC-Clothes, and DP3D. To assess the feasibility of CSCL in cases without clothing change, we also choose four SOTA short-term methods (i.e. PCB [37], HACNN [25], MGN [41], and Trans-ReID [13]) as competitors. The comparative results on the Market-1501 and DukeMTMC are presented in Table 3.

Based on the results in Table 2 and Table3, we have the following key observations: (1) In the cloth-changing setting, CSCL exceeds other competitors on PRCC, VC-Clothes, and DP3D by a large margin, achieving a rank-1 improvement of 4.9%/3.5%/9.6% and a mAP improvement of 6.8%/3.5%/10.9%. This is attributed to the powerful shape representation capability of the continuous correspondences. However, there is still a limitation to CSCL. Due to the poor quality of person images, the generated pseudo correspondences on LTCC are not reliable enough. Despite this limitation, CSCL still achieves comparable results with the SOTA method MVSE on LTCC, indicating a certain tolerance for vertex position errors. (2) CSCL generalizes well to the general ReID datasets where appearance features dominate, achieving comparable performance with the SOTA short-term methods. This is because the distribution of global RGB features is well preserved in the fusion stage.

Table 4: Ablation studies of different components in the CSCL framework. LNP represents linear projection, PE denotes positional embeddings, and SEN denotes the shape extraction network, respectively.

Models	CSE	SEN	CMF	Projection			PRCC		DP3D	
				LNP+PE	CP	LCP	Rank-1	mAP	Rank-1	mAP
1(Baseline)	X	-	-	-	-	-	32.8	37.1	19.2	13.0
2	✓	X	X	-	-	-	34.2	38.8	21.9	14.2
3	✓	✓	X	-	-	-	52.9	55.4	30.7	23.5
4	✓	✓	✓	✓	X	X	62.5	63.7	37.7	27.1
5	✓	✓	✓	✓	✓	✓	62.8	63.7	37.7	27.3
6	✓	✓	✓	X	X	✓	64.2	64.5	39.2	28.7

Table 5: Average Precision (AP) and Recall (AR) calculated at GPS thresholds ranging from 0.5 to 0.95 on multiple datasets.

Datasets	AP_{50}	AP_{75}	AP_{95}	AR_{50}	AR_{75}	AR_{95}
Market-1501	67.6	58.5	50.8	70.0	60.8	52.4
DukeMTMC	63.1	52.3	45.6	63.5	53.1	46.0
LTCC	59.2	49.8	39.5	60.3	51.7	39.8
PRCC	66.4	57.4	49.7	67.6	59.3	50.9
VC-Clothes	73.0	64.9	59.2	74.1	67.1	58.8
DP3D (w/o. \mathcal{L}_{cst})	84.0	74.6	65.8	84.9	76.0	66.2
DP3D	87.5	79.6	70.3	90.3	81.2	70.5

5.4 Ablation Studies

In this section, we carry out comprehensive experiments on PRCC and DP3D to validate: (1) the effectiveness of continuous surface embeddings, the cross-modality fusion module, and latent convolutional projection, which are abbreviated as CSE, CMF, and LCP respectively; (2) the influence of consistency loss on correspondence learning; (3) the impact of using different features for inference.

Effectiveness of CSE, CMF, and LCP. From Table 4, we observe that introducing continuous surface embeddings to the model with a proper shape extraction network (Baseline→Model3) remarkably improves the performance of the baseline model, with a rank-1/mAP improvement of 20.1%/18.3% on PRCC, and a rank-1/mAP

Table 6: Ablation studies of deploying different features for inference. CLS denotes the two learnable class tokens.

Features	PRCC		DP3D	
	Rank-1	mAP	Rank-1	mAP
RGB	55.9	57.7	36.8	26.3
CLS	61.1	61.8	37.9	27.4
Shape	42.5	45.9	23.9	16.8
RGB + Shape	61.4	62.6	37.4	27.2
CLS + Shape	63.5	64.2	38.7	28.4
RGB + CLS + Shape	64.2	64.5	39.2	28.7

improvement of 11.5%/10.5% on DP3D. This demonstrates that establishing pixel-wise and continuous correspondences complement rich and essential identity-related shape features for global RGB features. However, there is no significant improvement when directly downsampling the learned correspondences without a shape extraction network, and we will further analyze this issue in Section 5.5. Moreover, the cross-modality fusion module also brings significant improvement, which indicates that features of the two modalities become more compatible via cross-modality fusion. Furthermore, by comparing different feature projection methods for generating Q/K/V, we observe that LCP shows a certain degree of improvement over linear projection and convolutional projection. This is attributed to the inclusion of latent embeddings, which greatly facilitates the sharing among tokens.

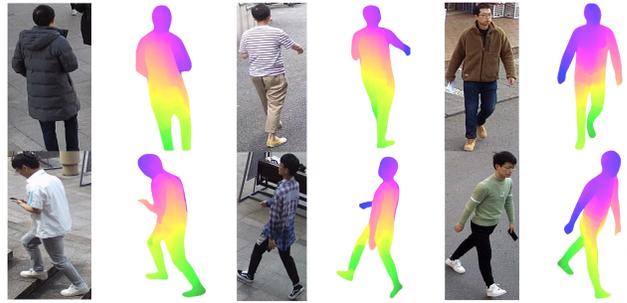
Additionally, we evaluate the quality of established correspondences on different ReID datasets in Table 5. By combining the results from Table 2 and Table 5, we can clearly observe a robust positive correlation between the quality of correspondences and the magnitude of performance improvement.

Influence of consistency loss. As shown in Table 5, the removal of consistency loss \mathcal{L}_{cst} from the correspondence learning process leads to a 5% decrease in vertex classification accuracy on DP3D, which indicates that performing consistency learning is beneficial for establishing reliable correspondences. From Table 2, we also observe that removing \mathcal{L}_{cst} results in a decline in the overall performance of ReID, verifying the importance of consistency learning for CSE.

Impact of using different features for inference. During inference, we select the model corresponding to Model 6 in Table 4 to verify the effectiveness of different features. As shown in Table 6, while relying solely on shape features is not reliable enough, the shape features can enhance the performance of other features. Concatenating global RGB features, shape features, and two class tokens results in the best performance at inference time.

5.5 Further Analysis

Visualization of Continuous Surface Embeddings. We employ PCA to reduce the dimension of continuous surface embeddings from $H \times W \times D$ to $H \times W \times 3$, where H and W denote the height and width of person images, D represents the embedding dimension. Visualization results on DP3D are presented in Figure 4. Since the color differences reflect the feature distances in the embedding space, we can clearly observe that the established 2D-3D correspondences between images pixels and the entire body surface are relatively smooth. Different from discretized UV mappings such as

**Figure 4: PCA visualization results of the learned continuous surface embeddings. The person images in each row are cross-appearance images of the same person in DP3D. We reduce the channel dimension of the learned continuous surface embeddings from 64 to 3 for visualization.**

the DensePose, the smooth and continuous 2D-3D correspondences can provide richer and more reliable global knowledge of human shape for cloth-changing ReID.

Identity modeling for shape features. Multi-modal auxiliary information itself is not sufficiently discriminative for the ReID task, making it necessary to conduct identity modeling. However, some existing CC-ReID methods, such as 3DSL, directly regulate multi-modal auxiliary features via ReID losses, which disrupts the distribution of shape space. As shown in Table 4, directly using downsampling operations without a proper shape extraction network (Model3→Model2) leads to significant performance degradation. We believe that multi-modal auxiliary features should first be mapped to an intermediary feature space before identity modeling to alleviate the incompatibility between feature spaces of different tasks, which is beneficial for the fusion of shape and global RGB features.

Future works. Current 3D shape-based ReID methods suffer from a huge domain gap between the RGB image space and the 3D shape space. Our work essentially targets at bridging the gap between these two spaces. Therefore, future works can consider transforming the surface embeddings into different forms of 3D shape features and assess their potential benefits for CC-ReID.

6 CONCLUSION

We have proposed a new shape embedding paradigm that establishes pixel-wise and continuous surface correspondences to mine fine-grained shape features for cloth-changing ReID. Moreover, an optimized cross-modality fusion module is designed to adaptively integrate shape features with global RGB features. To facilitate the research, we have constructed 3D Dense Persons (DP3D), which is the first cloth-changing ReID dataset with densely annotated 2D-3D correspondences and corresponding 3D meshes. Experiments on both cloth-changing and cloth-consistent ReID benchmarks demonstrate the robustness and superiority of our method.

ACKNOWLEDGMENTS

This work was supported in part by the Research Project of ZJU-League Research & Development Center, Zhejiang Lab under Grant 2019KD0AB01.

REFERENCES

- [1] Rıza Alp Güler, Natalia Neverova, and Kokkinos Iasonas. 2018. DensePose: Dense Human Pose Estimation in the Wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7297–7306.
- [2] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. 2021. Realtime Multi-Person 2D Pose Estimation Using Part Affinity Fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7291–7299.
- [3] Patrick P. K. Chan, Haorui Song, Peng Peng, Keke Chen, and Daniel S.Yeung. 2023. Learning Disentangled Features for Person Re-Identification under Clothes Changing. *ACM Transactions on Multimedia Computing, Communications, and Applications*. DOI:10.1145/3584359 (2023).
- [4] Jiaying Chen, Xinyang Jiang, Fudong Wang, Jun Zhang, Feng Zheng, Xing Sun, and Wei-Shi Zheng. 2021. Learning 3D Shape Feature for Texture-Insensitive Person Re-Identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8146–8155.
- [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. 2009. Imagenet: A Large-Scale Hierarchical Image Database. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- [6] Michela Farenzena, Loris Bazzani, Alessandro Perina, Vittorio Murino, and Marco Cristani. 2010. Person Re-Identification by Symmetry-Driven Accumulation of Local Features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2360–2367.
- [7] Zan Gao, Hongwei Wei, Weili Guan, Weizhi Nei, Meng Liu, and Meng Wang. 2022. Multigranular Visual-Semantic Embedding for Cloth-Changing Person Re-identification. In *Proceedings of the 30th ACM International Conference on Multimedia*. 3703–3711.
- [8] Ke Gong, Yiming Gao, Xiaodan Liang, Xiaohui Shen, Meng Wang, and Liang Lin. 2019. Graphonomy: Universal Human Parsing via Graph Transfer Learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 7450–7459.
- [9] Xinqian Gu, Hong Chang, Bingpeng Ma, Shutao Bai, Shiguang Shan, and Xilin Chen. 2022. Clothes-Changing Person Re-identification with RGB Modality Only. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1060–1069.
- [10] Rasmus Laurvig Haugaard and Anders Glent Buch. 2022. SurfEmb: Dense and Continuous Correspondence Distributions for Object Pose Estimation with Learnt Surface Embeddings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6749–6758.
- [11] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. 2017. Mask R-CNN. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2961–2969.
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 770–778.
- [13] Shuting He, Hao Luo, Pichao Wang, Fang Wang, Hao Li, and Wei Jiang. 2021. TransReID: Transformer-Based Object Re-Identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 15013–15022.
- [14] Alexander Hermans, Lucas Beyer, and Bastian Leibe. 2017. In Defense of the Triplet Loss for Person Re-Identification. arXiv preprint arXiv:1703.07737 (2017).
- [15] Peixian Hong, Tao Wu, Ancong Wu, Xintong Han, and Wei-Shi Zheng. 2021. Fine-Grained Shape-Appearance Mutual Learning for Cloth-Changing Person Re-Identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10513–10522.
- [16] Yan Huang, Jingsong Xu, Qiang Wu, Yi Zhong, Peng Zhang, and zhaoxiang Zhang. 2019. Beyond Scalar Neuron: Adopting Vector-Neuron Capsules for Long-Term Person Re-Identification. *IEEE Transactions on Circuits and Systems for Video Technology (TCSVT)* 30, 10 (2019), 3459–3471.
- [17] Yukun Huang, Zheng-Jun Zha, Xueyang Fu, and Wei Zhang. 2019. Illumination-Invariant Person Re-Identification. In *Proceedings of the 27th ACM International Conference on Multimedia*. 365–373.
- [18] Anastasia Ianina, Nikolaos Sarafianos, Yuanlu Xu, Ignacio Rocco, and Tony Tung. 2022. BodyMap: Learning Full-Body Dense Correspondence Map. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13286–13295.
- [19] Xuemei Jia, Xian Zhong, Mang Ye, Wenxuan Liu, Xenxin Huang, and Shilei Zhao. 2022. Patching Your Clothes: Semantic-Aware Learning for Cloth-Changed Person Re-Identification. In *International Conference on Multimedia Modeling (MMM)*. 121–133.
- [20] Bingliang Jiao, Lingqiao Liu, Liying Gao, Guosheng Lin, Ruiqi Wu, Shizhou Zhang, Peng Wang, and Yanning Zhang. 2022. Generalizable Person Re-Identification via Viewpoint Alignment and Fusion. arXiv preprint arXiv:2212.02398 (2022).
- [21] Xin Jin, Tianyu He, Kecheng Zheng, Zhiheng Yin, Xu Shen, Zhen Huang, Ruoyu Feng, Jianqiang Huang, Zhibo Chen, and Xian-Sheng Hua. 2022. Cloth-Changing Person Re-identification from A Single Image with Gait Prediction and Regularization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 14278–14287.
- [22] Hui Li, Yinglin Zheng, Zhaodong Tan, and Wenjin Deng. 2022. Improving Person Re-identification with Semantically Aligned Appearance Transformer. In *International Joint Conference on Neural Networks*.
- [23] Shuzhao Li, Huimin Yu, and Roland Hu. 2020. Attributes-aided Part Detection and Refinement for Person Re-Identification. *Pattern Recognition* 97 (2020), 4326–4335.
- [24] Wei Li, Rui Zhao, Tong Xiao, and Xiaogang Wang. 2014. DeepReID: Deep Filter Pairing Neural Network for Person Re-Identification. 152–159.
- [25] Wei Li, Xiatian Zhu, and Shangang Gong. 2018. Harmonious Attention Network for Person Re-Identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2285–2294.
- [26] Mengmeng Liu, Zhi Ma, Tao Li, Yanfeng Jiang, and Kai Wang. 2022. Long-Term Person Re-identification with Dramatic Appearance Change: Algorithm and Benchmark. In *Proceedings of the 30th ACM International Conference on Multimedia*. 6406–6415.
- [27] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. 2015. SMPL: A Skinned Multi-Person Linear Model. *ACM transactions on graphics (TOG)* 34, 6 (2015).
- [28] Zhangqiang Ming, Min Zhu, Xiangkun Wang, Jiaming Zhu, Cheng Junlong, Chengrui Gao, Yong Yang, and Xiaoyong Wei. 2022. Deep Learning-Based Person Re-Identification Methods: A Survey and Outlook of Recent Work. *Image and Vision Computing* 119:104394 (2022).
- [29] Zhangqiang Ming, Min Zhu, Xiangkun Wang, Jiaming Zhu, Cheng Junlong, Chengrui Gao, Yong Yang, and Xiaoyong Wei. 2022. IRANet: Identity-Relevance Aware Representation for Cloth-Changing Person Re-Identification. *Image and Vision Computing* 117:104335 (2022).
- [30] Natalia Neverova, David Novotny, Marc Szafrańiec, Vasil Khalidov, Patrick Labatut, and Andrea Vedaldi. 2020. Continuous Surface Embeddings. In *Proceedings of Neural Information Processing Systems (NeurIPS)*. 17258–17270.
- [31] Diederik P Kingma and Jimmy Ba. 2014. Adam: A Method for Stochastic Optimization. arXiv preprint arXiv:1402.6980 (2014).
- [32] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. 2019. Expressive Body Capture: 3D Hands, Face, and Body from a Single Image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10975–10985.
- [33] Xuelin Qian, Wenxuan Wang, Li Zhang, Fangrui Zhu, Yanwei Fu, Tao Xiang, Yu-Gang Jiang, and Xiangyang Xue. 2020. Long-Term Cloth-Changing Person Re-Identification. In *Proceedings of the Asian Conference on Computer Vision (ACCV)*. 71–88.
- [34] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi. 2016. Performance Measures and a Data Set for Multi-target, Multi-camera Tracking. In *Proceedings of the European Conference on Computer Vision*. Springer, 17–35.
- [35] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-Net: Convolutional Networks for Biomedical Image Segmentation. In *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. 234–241.
- [36] Xiujun Shu, Ge Li, Xiao Wang, Weijian Ruan, and Qi Tian. 2021. Semantic-Guided Pixel Sampling for Cloth-Changing Person Re-Identification. *IEEE Signal Process. Lett.* 28 (2021), 1365–1369.
- [37] Yifan Sun, Liang Zheng, Yi Yang, Qi Tian, and Shengjin Wang. 2018. Beyond Part Models: Person Retrieval with Refined Part Pooling (and A Strong Convolutional Baseline). In *Proceedings of the European Conference on Computer Vision*. Springer, 480–496.
- [38] Feitong Tan, Danhang Tang, Mingsong Dou, Kaiwen Guo, Rohit Pandey, Cem Keskin, Ruofei Du, Deqing Sun, Sofien Bouaziz, Sean Fanello, Ping Tan, and Yinda Zhang. 2021. HumanGPS: Geodesic PreServing Feature for Dense Human Correspondences. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1820–1830.
- [39] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All You Need. In *Proceedings of Neural Information Processing Systems (NeurIPS)*. 5998–6008.
- [40] Fangbin Wan, Yang Wu, Xuelin Qian, Yixiong Chen, and Yanwei Fu. 2020. When Person Re-Identification Meets Changing Clothes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 830–831.
- [41] Guanshuo Wang, Yufeng Yuan, Xiong Chen, Jiwei Li, and Xi Zhou. 2018. Learning Discriminative Features with Multiple Granularities for Person Re-Identification. In *Proceedings of the 26th ACM international conference on Multimedia*. 274–282.
- [42] Kai Wang, Zhi Ma, Shiyuan Chen, Jimni Yang, Keke Zhou, and Tao Li. 2020. A Benchmark for Clothes Variation in Person Re-Identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35, 12 (2020), 1881–1898.
- [43] Qizao Wang, Xuelin Qian, Yanwei Fu, and Xiangyang Xue. 2022. Co-Attention Aligned Mutual Cross-Attention for Cloth-Changing Person Re-Identification. In *Proceedings of the Asian Conference on Computer Vision (ACCV)*. 2270–2288.
- [44] Longhui Wei, Shiliang Zhang, Wen Gao, and Qi Tian. 2021. Person Transfer GAN to Bridge Domain Gap for Person Re-Identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 79–88.
- [45] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. 2021. CvT: Introducing Convolutions to Vision Transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 22–31.

- [46] Yuqiao Xian, Jinrui Yang, Fufu Yu, Jun Zhang, and Xing Sun. 2023. Graph-Based Self-Learning for Robust Person Re-Identification. In *IEEE Workshop on Applications of Computer Vision (WACV)*. 4789–4798.
- [47] Wanlu Xu, Hong Liu, Wei Shi, Ziling Miao, and Feihu Chen. 2021. Adversarial Feature Disentanglement for Long-Term Person Re-Identification. In *Proceedings of the 30th International Joint Conference on Artificial Intelligence*. 1201–1207.
- [48] Yuming Yan, Huimin Yu, Shuzhao Li, Zhaohui Lu, Jianfeng He, Haozhuo Zhang, and Runfa Wang. 2022. Weakening the Influence of Clothing: Universal Clothing Attribute Disentanglement for Person Re-Identification. In *Proceedings of the 31th International Joint Conference on Artificial Intelligence*. 1523–1529.
- [49] Qize Yang, Ancong Wu, and Wei-Shi Zheng. 2021. Person Re-Identification by Contour Sketch Under Moderate Clothing Change. *IEEE Trans. Pattern Anal. Mach. Intell* 43, 6 (2021), 2029–2046.
- [50] Zhengwei Yang, Xian Zhong, Zhun Zhong, Hong Liu, Zheng Wang, and Shin’ichi Satoh. 2023. Win-Win by Competition: Auxiliary-Free Cloth-Changing Person Re-Identification. *IEEE Transactions on Image Processing* 32 (2023), 2985–2999.
- [51] Hong-Xing Yu, Wu Ancong, and Zheng Wei-Shi. 2018. Unsupervised Person Re-Identification by Deep Asymmetric Metric Embedding. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 42, 4 (2018), 956–973.
- [52] Shijie Yu, Shihua Li, Dapeng Chen, Rui Zhao, Junjie Yan, and Yu Qiao. 2020. COCAS: A Large-Scale Clothes Changing Person Dataset for Re-Identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3400–3409.
- [53] Renjie Zhang, Yu Fang, Huaxin Song, Fangbin Wan, Yanwei Fu, Hirokazu Kato, and Yang Wu. 2023. Specialized Re-Ranking: A Novel Retrieval-Verification Framework for Cloth Changing Person Re-Identification. *Pattern Recognition* 134 (2023).
- [54] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, and Zhibo Chen. 2019. Densely Semantically Aligned Person Re-Identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 667–676.
- [55] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. 2015. Scalable Person Re-Identification: A Benchmark. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1116–1124.
- [56] Zhedong Zheng, Liang Zheng, and Yi Yang. 2019. Pedestrian Alignment Network for Large-scale Person Re-identification. *IEEE Transactions on Circuits and Systems for Video Technology* 29, 10 (2019), 3037–3045.