

Transformer for Object Re-Identification: A Survey

Mang Ye, Shuoyi Chen, Chenyue Li, Wei-Shi Zheng, *IEEE Senior Member*
David Crandall, *IEEE Senior Member*, Bo Du, *IEEE Senior Member*

Abstract—Object Re-Identification (Re-ID) aims to identify and retrieve specific objects from varying viewpoints. For a prolonged period, this field has been predominantly driven by deep convolutional neural networks. In recent years, the Transformer has witnessed remarkable advancements in computer vision, prompting an increasing body of research to delve into the application of Transformer in Re-ID. This paper provides a comprehensive review and in-depth analysis of the Transformer-based Re-ID. In categorizing existing works into Image/Video-Based Re-ID, Re-ID with limited data/annotations, Cross-Modal Re-ID, and Special Re-ID Scenarios, we thoroughly elucidate the advantages demonstrated by the Transformer in addressing a multitude of challenges across these domains. Considering the trending unsupervised Re-ID, we propose a new Transformer baseline, UntransReID, achieving state-of-the-art performance on both single-/cross modal tasks. Besides, this survey also covers a wide range of Re-ID research objects, including progress in animal Re-ID. Given the diversity of species in animal Re-ID, we devise a standardized experimental benchmark and conduct extensive experiments to explore the applicability of Transformer for this task to facilitate future research. Finally, we discuss some important yet under-investigated open issues in the big foundation model era, we believe it will serve as a new handbook for researchers in this field.

Index Terms—Object Re-Identification, Transformer, Literature Survey, Deep Learning



1 INTRODUCTION

The object re-identification (Re-ID), which aims at matching the same object (person [1], vehicles [2], etc) across multiple different views, has received significant attention in computer vision for a long time [3], [4], [5], [6], [7]. The core goal is to re-identify the same individual across different camera views or retrieve specific objects from large-scale image or video databases [4], [8], [9]. The main challenges of this task arise from the images captured from multiple cameras, which involve complex backgrounds, occlusions, changing lighting conditions, and diverse perspectives. Extensive efforts in the Re-ID field have been dedicated to addressing these issues [3], [10], [11]. Over the years, research in the field of Re-ID, especially concerning persons and vehicles [12], [13] has experienced rapid development, achieving human-level performance on a wide range of public datasets. In fact, Re-ID encompasses a wide range of subjects, including different types of objects (*e. g.*, animals and buildings), where the challenges are totally different. To better align with practical application requirements, existing Re-ID research is progressively expanding its focus towards broader scenarios and addressing more complex challenges. This includes dealing with limited annotation [14], [15], [16], diverse data modalities [17], [18], generalization of unknown scenarios [19], [20], and various applications [21], [22], [23].

Benefiting from the development of deep learning, the

studies in the Re-ID field have been dominated by Convolutional Neural Networks (CNNs) for a long time [3], [8]. Nevertheless, the emergence of the Vision Transformer [28], [29] has changed this situation. Transformer is a network architecture dispensing with recurrence and convolutions that relies entirely on attention mechanisms to model the global dependencies between inputs and outputs. Transformer has been introduced into the field of computer vision as a breakthrough, demonstrating remarkable performance in various visual tasks, including Re-ID. Intuitively, some works try to directly replace CNNs with existing vision transformers [29], [30], [31] as the feature extractor, which significantly improves the accuracy of Re-ID [21], [32], [33], [34], [35]. Unlike CNNs, Vision Transformer (ViT) [29] architecture imposes little structural bias to guide representation learning, which allows for diverse learning strategy design and broad task applicability [36]. This allows the Transformer to be flexibly applied to various scenarios including unsupervised, multimodal, etc. in Re-ID. Furthermore, ViT treats an image as a sequence of patch tokens rather than processing images pixel-by-pixel, which allows inputs of various sizes to be accepted without additional adjustments. The tokenization feature exhibits strong compatibility for personalized design and offers flexibility in organizing information [37], [38], [39]. These advantages facilitate the seamless integration of Transformer with Re-ID-specific designs, and novel research ideas for Re-ID that leverage the unique properties of transformers continue to emerge [21], [40], [41], [42], [43]. In recent years, Transformer-based Re-ID research has continuously created new records in recognition accuracy, showing a trend of being significantly superior to CNN-based methods in many aspects, as shown in Fig. 1. Due to the rapid proliferation of transformer-based Re-ID models, it is becoming progressively challenging to

- M. Ye, SY. Chen, CY. Li and B. Du are with the National Engineering Research Center for Multimedia Software, School of Computer Science, Hubei Luoqia Laboratory, Wuhan University, Wuhan, China. E-mail: {yemang, chenshuoyi, chenyueli, dubo}@whu.edu.cn
- WS. Zheng is with the School of Data and Computer Science, Sun Yat-sen University, Guangzhou, China (e-mail: wszheng@ieee.org)
- D. Crandall is with the Luddy School of Informatics, Computing, and Engineering, Indiana University. (Email: djcran@indiana.edu)

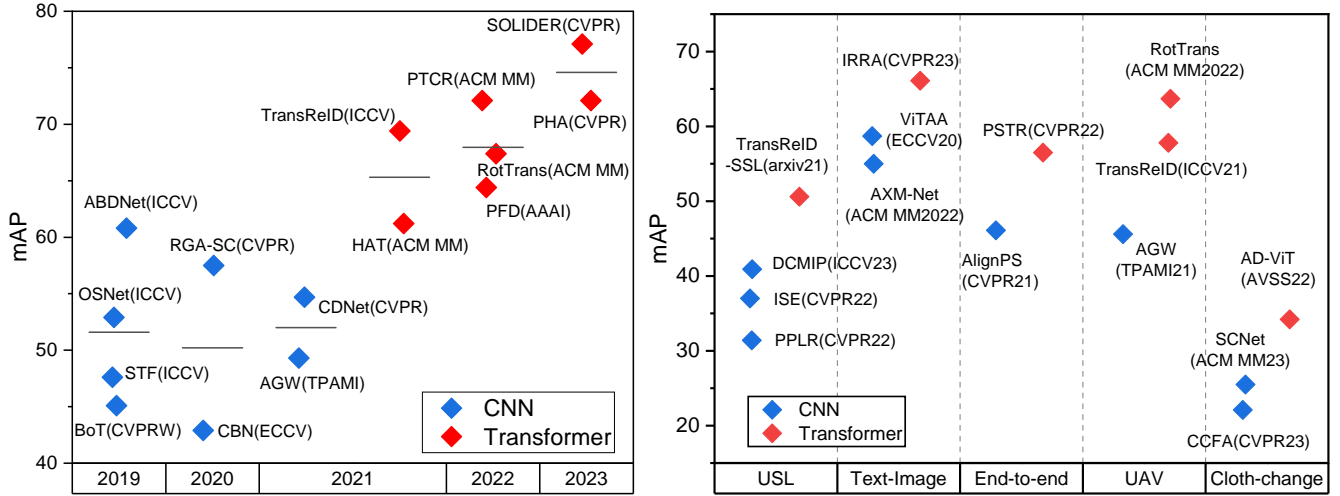


Fig. 1: (1) We show the performance of recent state-of-the-art methods on the widely-used person Re-ID dataset MSMT17 (left). The transformer-based methods have achieved a comprehensive lead in accuracy since 2021, while the CNN-based method for single-modal image Re-ID has not been studied. (2) We show state-of-the-art results of representative works in different Re-ID tasks: unsupervised (USL) Re-ID on MSMT17 [24] dataset, Text-Image on CUHK-PEDES [25], end-to-end person search on PRW [12], Re-ID in UAVs on PRAI-1581 [26], and cloth-changing Re-ID on LTCC [27].

stay abreast of the latest advancements. Consequently, there is an urgent need for a comprehensive survey of existing Transformer-based works, which would greatly benefit the community in the new era.

Existing Re-ID surveys [3], [8], [44], [45] predominantly focus on deep learning methods based on CNNs and tend to narrow their scope to specific objects, with a primary emphasis on persons or vehicles. On the contrary, this survey is mainly oriented to the application of emerging transformer technology in Re-ID and covers a wider range of objects (person, vehicle, and animal), which is more innovative and comprehensive. Recognizing the significant potential and promise demonstrated by numerous Transformer-based studies in various vision applications, we systematically organize and review the growing research works on Transformer for Re-ID in recent years to gain valuable insights. Differing from existing surveys, the distinctive features and primary contributions of our survey are as follows:

- We conduct an in-depth analysis of the strengths of Transformer and summarized the research efforts since its introduction into Re-ID field across four extensively studied Re-ID directions, including image/video-based Re-ID, Re-ID with limited data/annotations, cross-modal Re-ID, and special Re-ID scenarios. It demonstrates the success of Transformer in Re-ID and underscores its potential for future advancements.
- We propose a Transformer-based unsupervised baseline for trending unsupervised Re-ID since the Transformer has not been extensively explored in existing works. It achieves competitive performance on both single- and cross-modal unsupervised Re-ID tasks.
- We particularly delve into animal re-identification, an area that has received significantly less attention compared to persons and vehicles. It presents numerous challenges and unresolved issues. We develop unified experimental standards for animal Re-ID and evaluate the feasibility of employing Transformer in this context,

laying a solid foundation for future research. A periodically updated website will available at <https://github.com/JigglypuffStitch/Animal-Re-ID> for researchers in this field.

The rest of this survey is organized as follows: In § 2, we briefly review the development of the Re-ID field before Transformer era, while introducing Transformer in vision and providing a detailed analysis of its numerous advantages. A comprehensive analysis of Transformer in Re-ID and a powerful Transformer baseline for single/cross-modal unsupervised Re-ID are presented in § 3. The progress in Animal Re-ID and the evaluation of the applicability of Transformer to this task are introduced in § 5.

2 BACKGROUND

2.1 A Brief Review of Re-ID Before Transformers

This subsection first summarizes the basic definition and challenges of object Re-ID (§ 2.1.1). Then, we give a general review of the past Re-ID research works dominated by CNNs and discuss the limitations of CNN-based Re-ID methods in some aspects (§ 2.1.2).

2.1.1 Object Re-Identification

Definition. Given a query q , the goal of object re-identification is to retrieve specific objects from a gallery set $\mathcal{G} = \{g_i | i = 1, 2, \dots, N\}$ of N descriptors. An important feature of Re-ID is to identify the objects in different cameras with non-overlapping views. The query q can be an image, a video sequence, a text description, a sketch, or a combination of different forms, *etc.* [8], [17], [40], [46], [47], [48], [49]. The identity of the query q can be formulated as:

$$I = \arg \min_{g_i} \text{dis}(q, g_i), g_i \in \mathcal{G}, \quad (1)$$

where $\text{dis}(\cdot, \cdot)$ is an arbitrary distance metric.

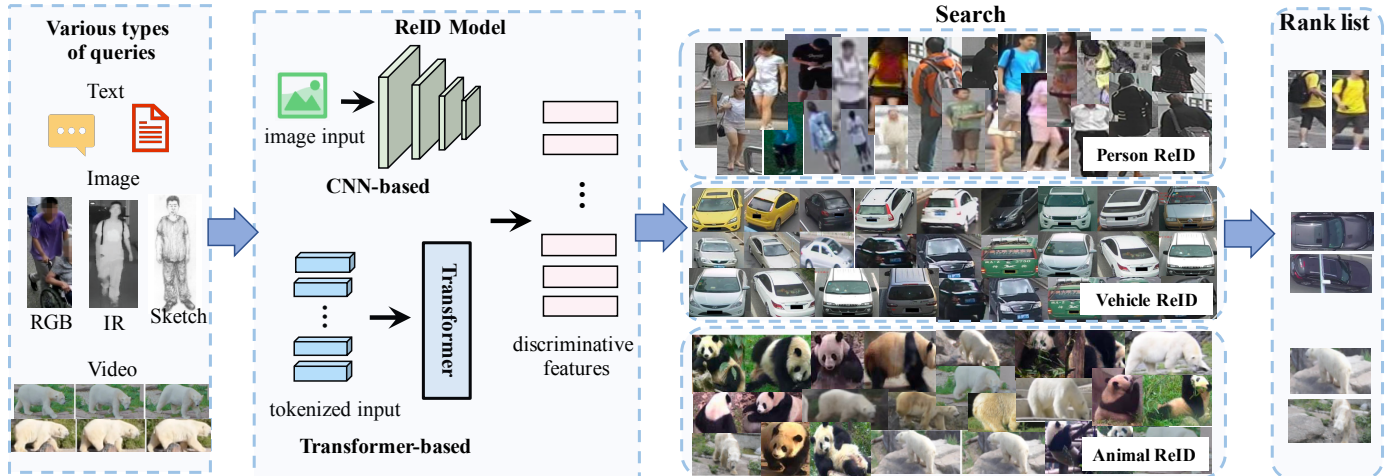


Fig. 2: **General object Re-ID process.** Given a query that can be any type of image, text, video, etc., the goal of Re-ID is to search for the specific object from gallery data collected by different cameras.

Challenges. The flow of a basic Re-ID system is shown in Fig. 2. The early Re-ID tasks mainly focused on the pre-detected object bounding boxes obtained from the original image or video and used the appearance information to match the corresponding individual. Due to the complex conditions of image acquisition, the difficulty of Re-ID in this period involves occlusion, illumination variation, resolution difference, camera view variation, and background clutter [5], [8], [44]. Furthermore, new challenges continue to emerge with the urgent growth of practical application requirements [16], [18], [50]. For example, with the widespread application of drone surveillance recently [13], [26], [51], discriminative information has been greatly reduced under extreme bird’s-eye view angles [52], [53]. In the data processing stage, for special cases where conventional visible light images are not available, valid object information might be represented by different modalities such as infrared images, text descriptions, sketches, and depth images [49], [54], [55]. The cross-modal gaps of object Re-ID in the modal heterogeneous scene lead to great intra-class differences across varying modalities [25], [56]. In addition, due to artificial cross-camera correlation of objects, high labeling costs [57], [58] and unavoidable noise problems [59] make it difficult to achieve large-scale expansion [60], [61], [62]. In the retrieval phase, the varying real environmental domains may cause the inapplicability of the Re-ID model [63], [64], [65] and long-term Re-ID suffers from appearance information changes [27], [66].

2.1.2 CNN-based Re-ID Methods.

Under the mainstream trend of deep learning, the object Re-ID steps are generalized as different steps, including data processing, model training and descriptor matching. Most existing methods take training a strong Re-ID model as the core goal. In fact, CNNs have dominated Re-ID studies for a long period. In this section, we focus on reviewing the progress of object Re-ID which is highly related to CNNs. Considering different application requirements, Ye *et al.* proposed to divide Re-ID technology into two subsets, closed-world and open-world [3].

Closed-world refers to supervised learning methods based on well-labeled visible images captured by common video surveillance [47]. With the aid of labels, many approaches model Re-ID as a classification, verification or metric learning problem, using CNNs (*i.e.* ResNet [67]) to learn discriminative feature representations from the training data [46], [68], [69]. On the basis, learning local features such as image slices [11], [70], semantic parsing [71], [72], pose estimation [73], [74], region of interest [75] and key points [76], [77] to further mine fine-grained information are typical ideas in Re-ID. At the CNN backbone level, some people try to directly improve the convolutional layer and residual block [78], and a large number of studies introduce attention modules in CNNs to capture the relationship between different convolutional channels, feature maps, and local regions [3], [79], [80], [81], [82]. On the other hand, for video sequence input with temporal information, the major limitation of CNNs is that it can only process spatial dimension information. Some video Re-ID works introduce RNN or LSTM for sequence modeling [83], [84], [85].

Open-world technologies usually target on more complex and difficult scenarios, including cross-modal Re-ID, unsupervised learning, domain generalization, etc. (1) *Cross-modal Re-ID.* In recent years, cross-modal Re-ID of visible-infrared [17], [86], [87] and text-image [88], [89], [90] has received more and more interests. For visible-infrared Re-ID, researchers design single-stream [91], dual-stream [92], [93], [94] and other different structures [95] based on CNN to learn modality-sharing and modality-specific feature representations. To reduce the difference between modalities, a class of methods implements modal conversion or style transformation through GAN [96], [97], [98] or special augmentation strategies [17], and then performs subsequent CNN-based feature representation learning. Text-to-image Re-ID mainly focuses on the cross-modal alignment module design based on the visual and text features extracted from each modality backbone [99], [100]. In addition, many works also introduce attention mechanisms to enhance local information matching [101], [102]. (2) *Unsupervised learning.* This approach alleviates the label insufficiency issue [103], which now has been a trending topic due to its benefits in

large-scale applications. It mainly includes two categories: unsupervised domain adaptation [104], [105], [106] and pure unsupervised learning [107], [108]. In addition to transferring knowledge from labeled source datasets to unlabeled target datasets [24], [109], most of existing methods learn feature representations purely from unlabeled images [15]. The core idea is to use the features extracted by CNNs to perform clustering to obtain pseudo-labels as label supervision, some of which focus on generating high-quality pseudo-labels [62], [110], [111], and others improve clustering algorithms and training strategies [61], [103], [112]. (3) *Other open scenes*. In recent years, an increasing number of research efforts have shifted towards open scenarios, such as cloth-changing Re-ID and domain-generalizable Re-ID. To facilitate the research, many new cloth-changing datasets [27], [113] have been introduced. The key to addressing the cloth-changing problem lies in learning clothing-agnostic features, and straightforward approaches involve augmenting the data by introducing a variety of clothing types [114], [115]. Many works try to utilize auxiliary information, such as human parsing [116], [117], gait [118], shape [119] to guide the CNN model to focus on identity-related features. Domain generalization is also highly aligned with practical application requirements. Some research endeavors focus on creating large-scale and diverse synthetic data [63], while others seek to enhance the generalization capabilities of CNN models through meta-learning [65], [120] or disentanglement techniques [19].

2.2 Understanding and Analysis of Transformer

The introduction of Vision Transformer opens novel directions for Re-ID studies, especially in challenging scenarios. In this subsection, we first give the basic concept of the Transformer (§ 2.2.1). In order to demonstrate the superiority of the Transformer, we provide a comprehensive comparison between Transformer and CNN and analyze it in terms of network architecture, modeling capabilities, scalability, flexibility, and special properties (§ 2.2.2).

2.2.1 Transformer Concepts

Original Transformer. The original Transformer [28] was proposed in the field of natural language processing (NLP), which is the first sequence transduction model based on the attention mechanism. It completely abandons the dominant sequence transduction models based on complex recurrent and convolutional neural networks and achieves new state-of-the-art levels in multiple NLP tasks. The transformer is essentially an encoder-decoder structure, in which both encoder and decoder are composed of multiple stacked transformer layers [28], [37], [38]. Each transformer layer consists of two sub-layers: a multi-head self-attention mechanism and a position-wise fully connected feed-forward network. Self-attention plays a crucial role in Transformer, which enables each element to learn to gather from other tokens in the sequence. Multi-head self-attention can create multiple attention matrices in a layer, and with multi-head, the self-attention layer will create multiple outputs to guarantee the diverse capability. The two sub-layers perform residual connection [67] for stability, followed by layer normalization. Transformer accepts tokenized sequences as

input. To make efficient use of sequence order, an optional positional encoding (relative or absolute) needs to be added. Transformers are used in machine translation tasks, where the encoder extracts features from input with positional encodings, and the decoder uses these features to produce output. Since Transformer was proposed, it has gradually become mainstream and most of the subsequent NLP research has been reprocessed on its basis [121].

Vision Transformer. The emergence of Vision Transformer (ViT) [29] is a significant breakthrough in the field of computer vision. Different from the previous work that embeds the attention module in the CNN, it applies the pure Transformer to the image patches with a simple idea and has achieved remarkable results. Specifically, given an image $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$, $H \times W$ and C represent the image resolution and the number of channels respectively. In order to adapt to the tokenized sequence input of the Transformer, ViT designs a patch embedding operation that divides an image into N patches, where the size of each patch is $P \times P$. These patches are projected into the D -dimensional space after linear transformation as the input of ViT, which is a sequence composed of N D -dimensional vectors, denoted as $\mathbf{x} \in \mathbb{R}^{N \times D}$. A special learnable embedding called class token is set for classification, which is directly concatenated to the patch embedding. Following a similar line of the original Transformer, the position embedding is also added to each patch embedding to preserve the spatial position information of the image which is represented as $E_{pos} \in \mathbb{R}^{(N+1) \times D}$. ViT adopts the same structure as the encoder of the original Transformer as a feature extractor. Following the ViT paradigm, a series of subsequent ViT variants are proposed for various vision tasks, leading to significant advancements [38].

2.2.2 Superiority of Vision Transformer

We provide a detailed analysis of the strengths of Transformer in the vision perspective to facilitate the subsequent elaboration of its robust performance in addressing complex and dynamic Re-ID scenarios.

Powerful Modeling Capabilities. Different from the standard CNNs, which is limited to the local receptive field, it is extremely difficult to establish long-distance relationships at the early stage. As a result, the performance of CNNs is limited for challenging scenarios. In contrast, the powerful modeling ability of Transformer is reflected in the local-global duality [36]. Specifically, the modeling of images mainly involves pixel level and object level in vision tasks with images or videos. The attention of Transformer is flexibly designed to process information from any image region and can model any relationship between pixels-pixels, objects-pixels and objects-objects. Moreover, CNN can only build hierarchical representations from local to global, whereas Transformer has the flexibility to integrate global information at any stage [30], [39]. For Re-ID, both global and local modeling are essential for learning discriminative features to distinguish high-similar inter-class objects.

Diverse Unsupervised Learning Paradigms. Due to the expensive and time-consuming nature of acquiring large amounts of high-quality annotated data, unsupervised learning can develop more generalized feature representations without relying on annotations. The great success

of Transformers in NLP has largely benefited from self-supervised learning, which provides a solid foundation for the self-supervised research in Vision Transformer. Unsupervised learning in the field of computer vision has primarily been centered around contrastive learning, and the introduction of Transformer has made it feasible to incorporate mainstream generative learning approaches from NLP, such as masked autoencoders [122]. Besides, discriminative self-supervised methods reveal some new characteristics in Transformer models, such as clear object boundaries [123]. In general, unsupervised learning with its cost-effectiveness and generalization capabilities, emerges as a future trend, and transformers hold a unique advantage within this trend.

Multi-modal Uniformity and Versatility. In practical applications, single modal data may lack information or be ambiguous. Multi-modal learning allows models to leverage diverse types of data, enabling them to capture rich and comprehensive information for a better understanding of complex real-world scenarios. Compared to CNN, which is primarily designed for processing image modalities, Transformer exhibits significant versatility across multiple modalities. Multi-modal information can be transformed into a token sequence or latent space features within the same semantic space and input into Transformer for encoding. Transformer can be considered as a fully connected graph, where each token embedding is represented as a node in the graph and the relationships between these embeddings can be described by edges. This property enables Transformer to function within a modality-agnostic pipeline that is compatible with various modalities [37]. Especially in the combination of vision and language, Transformer promotes many new ideas for solving cross-modal Re-ID challenges.

High Scalability and Generalization. With the continuous increase in data, the future demand for highly scalable models becomes increasingly urgent to adapt effectively to the growing scale of data. Moreover, the generalization capability is crucial for stable performance in dynamic and unknown environments. Numerous recent studies have shown that the powerful scalability of Transformer in terms of large models and big data has achieved incredible results [124], [125]. Zhai *et al.* [126] successfully trained a Vision Transformer model with 2 billion parameters, achieving a new record of 90.45% top-1 accuracy on ImageNet. Transformers hold immense potential for larger and more versatile models. The encoder-decoder structure of the Transformer, coupled with the joint learning of decoder embeddings and positional encoding, can seamlessly unify various tasks. Additionally, its powerful cross-modal learning capabilities offer further possibilities for expanding Re-ID applications.

3 TRANSFORMER IN OBJECT RE-ID

In this section, we comprehensively review the latest research on Transformer-based Re-ID. Considering different types of challenges and the application of Transformer’s diverse advantages, we introduce Re-ID methods divided into four scenarios: regular images or videos with annotations (§ 3.1), limited data or limited annotations (§ 3.2), multimodal data (§ 3.3), and special settings (§ 3.4).

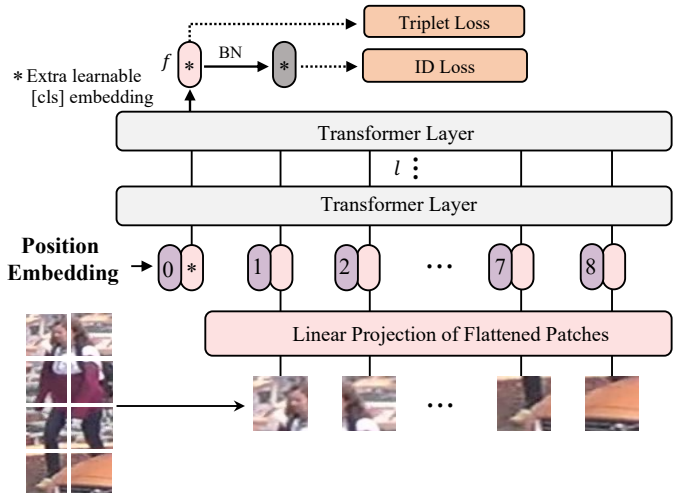


Fig. 3: The first pure transformer baseline for object Re-ID [32]. The Vision Transformer backbone [29] is adopted as a feature extractor, optimized with ID loss and triplet loss [6] widely used in Re-ID.

3.1 Transformer in Image/Video Based Re-ID

In this subsection, we summarize the progress of transformers under the general supervised setting of image-based (§ 3.1.1) and video-based (§ 3.1.2) Re-ID. For image-based Transformer Re-ID methods, we first review different structural designs at the backbone level for discriminative Re-ID feature extraction. In addition, the tokenized embeddings and attention mechanism in Transformer provide strong flexibility for representation learning. We comprehensively summarize the methods of exploiting Transformer properties for Re-ID-specific design. In video Re-ID, Transformer-based methods have shown great superiority over CNNs on modeling the spatio-temporal cues.

3.1.1 Transformer in image-based supervised Re-ID

Architecture Improvements. A number of recent studies have shown that applying Vision transformer as a feature extractor in Re-ID can achieve high accuracy [32], [33], [34], [129]. Many recent Re-ID methods are dedicated to designing special Transformer architectures to build stronger backbones [43], [127], [128]. The first to introduce Vision Transformer in the Re-ID field is TransReID [32], which preserves two advantages of Transformer at the architectural level. Compared with CNNs, the multi-head self-attention scheme in Transformer captures long-distance dependencies so that different object/body parts can be better focused. In addition, Transformer retains more detailed information without down-sampling operators. Therefore, it builds a pure Transformer baseline for supervised image-based single modality object Re-ID, following in a similar way to ViT [29], as shown in Fig. 3. Even simply replacing the feature extraction network in basic Re-ID methods with vision transformers, the performance on multiple vehicle and person Re-ID datasets is comparable to state-of-the-art methods, reflecting the strong potential of Transformers for Re-ID tasks. Inspired by this, some subsequent methods design special transformer architectures such as pyramid structure [43], hierarchical aggregation [127], [130], graph

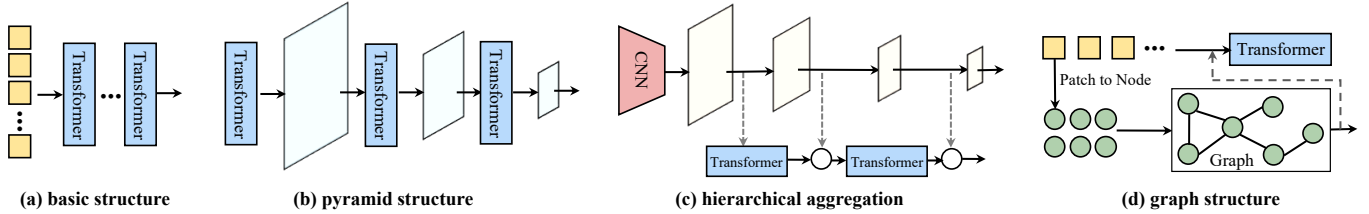


Fig. 4: **Different Transformer architectures designed for image-based Re-ID.** (a) The basic Re-ID baseline based on Vision Transformer [32]. (b) Pyramid Transformer for learning multi-scale features [43]. (c) Transformer and CNN hybrid architecture for aggregating hierarchical features [127]. (d) Combination of graph structure and Transformer [128].

structure [128], *etc.*, while some other methods intent to improve the attention mechanisms [131], [132], [133], [134].

Considering the mutual cooperation of CNN and Transformer, Li *et al.* [43] develop a pyramidal transformer structure like CNN to learn multi-scale features and improve the patch embedding process, utilizing convolution with the anti-aliasing block to capture translation-invariant information. Similarly, based on hierarchical features extracted by CNN, HAT [127] is proposed to aggregate features of different scales in a global view with the help of Transformer. GiT [128] introduces graphs in transformers to mine relationships of nodes within the patch.

For the improvement of attention schemes in Transformer, Zhu *et al.* [132] present a dual cross-attention learning strategy by emphasizing the interaction between the global image and local high-response regions and the interaction between image pairs. Considering the impact of variable appearance of the same identity, Shen *et al.* [134] introduce cross-attention in the transformer encoder to merge information from different instances. From the perspective of improving Transformer efficiency in Re-ID, Tian *et al.* [133] present a hierarchical walking attention, by introducing a prior as an indicator to decide whether to skip or calculate a region’s attention matrix in the image patch. For lightweight Re-ID models, Mao *et al.* [135] design an attention map-guided Transformer pruning method, so that the models can be deployed on edge devices with limited resources. It removes redundant tokens and heads in a hardware-friendly manner, achieving the goal of reducing the inference complexity and model size without sacrificing the accuracy of Re-ID.

Re-ID-specific Design. For different objects (especially person) and Re-ID-specific challenges (occlusions), many works explore the application of Transformer to make Re-ID-specific adaptations [136], [137], [138]. For the most crucial local discriminative information mining for Re-ID, Vision Transformer naturally has attention and patch embeddings, which can be easily used to capture local discriminative information to enhance the representation [32], [136], [137], [138]. Furthermore, the disentanglement of some key information can be modeled by the encoder-decoder structure in transformer [138], [139], [140]. Besides, the structure prior [141] or task specialties [21] of different objects is also important for Transformer design.

The effectiveness of learning local feature representation has been proven by extensive Re-ID research, and adding an attention mechanism to the CNN to focus on more discriminative information is a mainstream practice

in previous studies (§ 2.1.2). However, many new ideas that use the special properties of Transformers to learn local features are also emerging and growing rapidly [142]. TransReID [32] designs shift and patch shuffling operations on the transformer baseline to learn local features, which is conducive to enhancing disturbance invariance and robustness. Zhu *et al.* developed an Auto-Aligned Transformer (AAformer) [136] to adaptively locate the human parts. It learns local representations by introducing learnable part tokens to the transformer and integrates part alignment into self-attention. Zhang *et al.* revealed that self-attention leads to the inevitable dilution of high-frequency components of images. To enhance the feature representation of high-frequency components that are important to Re-ID, they proposed to use Discrete Haar Wavelet Transform (DHWT) [143] to split the patches of high-frequency components as the auxiliary information [144].

Recently Transformer-based methods have made great contributions to address the occlusion challenge of Re-ID [135], [139], [140], [145], [146], [147], [148], [149]. Extracting partial region representation features is also a good solution to the challenge of object occlusion in Re-ID. Part-Aware Transformer (PAT) [138] is proposed to exploit the Transformer encoder-decoder architecture to capture different discriminative body parts, where the encoder is used to obtain pixel context-aware feature maps and the decoder is used to generate part-aware masks. For the widely studied person object, due to the occlusion noise or the occlusion region is similar to the target, an intuitive solution is to guide local feature learning with the help of human pose information. Wang *et al.* [139] also adopt the transformer encoder-decoder structure to present a pose-guided feature disentangling method. With the keypoint information captured by the pose estimator, a set of learnable semantic views are introduced into the decoder to implicitly enhance the disentangled body part features. Similarly, assisted by motion information, Zhou *et al.* [140] utilize keypoint detection and part segmentation for Transformer encoder-decoder modeling. Besides, some transformer studies analyze occlusion problems from other perspectives [146], [149]. Xu *et al.* [146] design a feature recovery Transformer (FRT) to recover occluded features using nearby target information. Considering the similarity between the local information of each semantic feature in k-nearest neighbors and the query, FRT filters out noise to restore the occluded query feature. Cheng *et al.* [149] utilize knowledge learned from different source datasets to generate reliable semantic clues to alleviate domain differences between off-the-shelf semantic

models and Re-ID data. Transformer allows human parsing results to be embedded as learnable tokens into the input, where a weighted sum operation is employed to integrate parsed information from multiple sources.

3.1.2 Transformer in Video-based Re-ID

Transformer for Post-processing. Video Re-ID aims to fully exploit the temporal and spatial interactions of frame sequences to extract more discriminative representations [154]. Compared with CNN-based methods that require additional models to encode time information, Transformer is proposed as a powerful architecture for processing sequence data, which has inherent advantages. The global attention mechanism in Transformer can be easily adapted to video data to capture spatio-temporal dependencies [153]. A group of Transformer-based video Re-ID methods are hybrid architectures [151], [155], [156]. They are mainly based on the features initially extracted by CNN, using Transformer to model the long-range context of the video as post-processing. Zhang *et al.* [156] designed a two-stage spatio-temporal transformer module for patch tokens converted from CNN feature maps, where the spatial transformer focuses on object regions with different backgrounds, while the subsequent temporal transformer focuses on video sequences to exclude noisy frames. Also based on the features extracted by CNN, Liu *et al.* [155] present a multi-stream Transformer architecture that emphasizes three perspectives of video features via spatial Transformer, temporal Transformer and spatio-temporal Transformer. A cross-attention based strategy is designed to fuse multi-view cues to obtain enhanced features. Additionally, some studies consider the complementary learning of CNN and Transformer in space and time. Specifically, DCCT [152] introduces self-attention and cross-attention to the features extracted by two separate networks to establish a spatial complementary relationship, and designs hierarchical aggregation based on a temporal Transformer to integrate two temporal features. DenseIL [151] is a hybrid architecture consisting of a CNN encoder and a Transformer decoder with dense interaction, where the CNN encoder extracts discriminative spatial features while the decoder aims to densely model spatio-temporal interactions across frames.

Pure Video Transformer. The hybrid architecture is difficult to overcome the intrinsic limitations of CNN for perceiving long-distance information. Some recent work attempts to explore the application of pure Transformer architecture to video Re-ID [153], [154]. Tang *et al.* [153] design a multi-stage Transformer framework by taking advantage of Vision Transformer’s class token to facilitate the aggregation of various information. At different stages, the learning of attribute-associated information, identity-associated information and attribute-identity-associated information is guided respectively. Different from the mainstream divide-and-conquer strategy that tackles feature representation and feature aggregation separately that fail to simultaneously solve temporal dependence, attention and spatial misalignment, a contextual alignment Vision Transformer (CAViT) [154] is proposed for spatial-temporal joint modeling. To jointly model spatio-temporal cues, it replaces self-attention with temporal-shift attention based on a pure Transformer architecture to align objects in adjacent frames.

3.2 Transformer in ReID with Limited Data/Annotations

In this survey, limited annotation usually corresponds to unsupervised learning Re-ID technology (§ 3.2.1), while limited data mainly focuses on domain generalization in Re-ID (§ 3.2.2). In fact, Transformer is still in the preliminary exploration in such Re-ID scenarios, with a small number of works but showing great potential.

3.2.1 Transformer in Unsupervised Re-ID

Self-supervised Pre-training. Generally, the existing unsupervised Re-ID methods mainly rely on the features extracted by CNN to cluster and generate pseudo-labels as label supervision, in which CNN is supervisedly pre-trained on ImageNet [15], [112]. However, supervised pre-training focuses on coarse category-level distinction, which reduces the rich visual fine-grained information in images. This fine-grained cues are crucial for Re-ID tasks with large intra-class variation. A class of studies of Transformer in unsupervised Re-ID emphasizes self-supervised pre-training to obtain a better initialization model and reduce the domain difference between ImageNet data and Re-ID data [41], [159]. The success of self-supervised Transformer in vision tasks provides a lot of modeling and training experience for unsupervised Re-ID research [38]. The advantages of Vision Transformer in unsupervised learning are reflected in several aspects: 1) The strong scalability of the Transformer model for large-scale unlabeled data. The self-supervised learning can fully use the representation ability of the models with Transformer architecture [123]. 2) The flexibility of the Transformer structure provides more diverse self-supervised paradigms, which are extremely challenging for CNNs to complete [122].

With the emergence of LUPerson [58], a large-scale unlabeled dataset specifically for person Re-ID, Luo *et al.* [41] began to initially explore effective Transformer self-supervised pre-training paradigms for Re-ID, achieving significant results. They first conduct extensive experiments to investigate the performance of CNNs and Transformer using different Self-Supervised Learning (SSL) methods on ImageNet and LUPerson pre-training datasets. In addition, they promote the stability and domain invariance of Transformer by designing the IBN-based convolution stem to replace the standard patchify stem in ViT to enhance the local feature learning. The conclusion is that Transformer is ahead of CNN in terms of pre-training. Notably, under the fully unsupervised condition of using DINO [123] to pre-train Transformer on LUPerson and fine-tuning with common unsupervised Re-ID method [112], the performance of Re-ID is even competitive with the state-of-the-art supervised Re-ID method. It can be regarded as a major breakthrough in the field of unsupervised Re-ID. On this basis, the later PASS [159] further integrated Re-ID-specific part-aware properties in the self-supervised Transformer pre-training. Inspired by DINO [123] which develops a simple strategy for label-free self-distillation, PASS introduces several learnable tokens to extract part-level features, further reinforcing fine-grained learning for Re-ID. Specifically, it divides the image into several fixed overlapping local regions and randomly crops local views from them while the global view is randomly cropped from the whole image. In knowledge distillation,

TABLE 1: Representative Transformer methods based on image/video Re-ID.

Category	Focal Point	Object	Transformer Type	Method	Publication
Image-based Re-ID					
Architecture Design	Pure Transformer Re-ID baseline	Vehicle&Person	ViT	transReID [32]	ICCV 2021
	Hierarchical feature aggregation	Person	Hybrid	HAT [127]	ACMMM 2021
	Introducing the benefits of CNN	Person	PVT	PTCR [43]	ACMMM 2022
	Improvements to attention	Vehicle&Person	ViT,DeiT	DCAL [132]	CVPR 2022
	Integrating graph structure	Vehicle	ViT	GiT [128]	TIP 2023
Re-ID-specific Design	Partial representation learning	Person	Decoder	PAT [138]	CVPR 2021
	Introducing auxiliary information (e.g. pose, keypoint, etc)	Person	Encoder-Decoder	PFD [139]	AAAI 2022
	Explicit modeling of relationships between persons	Person	Transformer	NFormer [150]	CVPR 2022
	Learning rotation-invariant features	Person&Vehicle	ViT	RotTrans [21]	ACMMM 2022
	High-frequency augmentation	Person	ViT	PHA [144]	CVPR 2023
Video-based Re-ID					
Combination of CNN & Transformer	Transformer for post-processing spatio-temporal interaction	Person	Decoder	DenseLL [151]	ICCV 2021
	Coupled CNN-Transformer for complementary learning	Person	ViT/Swin/DeiT	DCCT [152]	TNNLS 2023
Pure Transformer	Spatial-temporal aggregation	Person	ViT	MSTAT [153]	TMM 2022
	Spatial-temporal joint modeling	Person	ViT	CAViT [154]	ECCV 2022

all views are passed through the student and only the global view is passed through the teacher.

The pre-trained Transformer serves as a powerful initialization model compared with previously widely-used ImageNet pretraining. It can be fine-tuned with different Re-ID supervised learning or unsupervised learning methods in downstream tasks. As shown in Tab. 2, the performances of corresponding methods have been greatly improved. We believe these research will be an advancement for the Re-ID community, allowing future work to be performed with better pre-trained models.

Unsupervised Domain Adaptation. Transformer has received limited attention for another widely studied unsupervised domain adaptation (UDA) problem in unsupervised Re-ID, with a small amount of work on vehicles and persons respectively [160], [161]. Different from the previous Re-ID method based on domain alignment to guide feature learning to achieve distribution consistency at the domain level or identity level, Wang *et al.* [160] oriented person Re-ID to achieve fine-grained domain alignment between different body parts with the help of Transformer. They embed the transformer layer into the feature extraction backbone and discriminators respectively, where the backbone obtains the class tokens representing each body part and the discriminators extract the domain information contained in each body part. Dual adversarial learning is introduced in the backbone and discriminator to align each class token of a target domain sample with the corresponding class token in the source domain. Another vehicle-oriented Transformer work belongs to the clustering-based UDA solution [161]. The core idea is to make the Transformer adaptively focus on the discriminative part of the vehicle in each domain through a joint training strategy. To achieve dynamic knowledge transfer, both source and target images are simultaneously fed into a shared CNN to obtain feature

maps, and the Transformer encoder-decoder architecture is subsequently introduced to generate a global feature representation integrating contextual information from the feature maps. Based on Transformer, a learnable domain encoding module similar to positional encoding is added to better utilize the specific characteristics of each domain.

3.2.2 Transformer in Generalized Re-ID

The application of Transformer promotes new ideas of Re-ID in the challenging problem of domain generalization (DG) [162]. Completely different from mainstream research that using Transformer for feature representation learning, TransMatcher [162] studies Transformer for image matching and metric learning for a given image pair from the perspective of generalizability. For Re-ID, a typical image matching and metric learning problem, the Transformer encoder can only facilitate the feature interaction between different positions within an image but fail to realize the interaction between different images. Liao *et al.* citeliao2021transmatcher also demonstrate that Directly applying vanilla vision Transformer ViT through classification training pipeline, it will result in poor generalization to different datasets. Inspired by the cross-attention module in the Transformer decoder that enables cross-interaction between query and encoded memory, they attempt to use actual image queries instead of learnable query embeddings as the input to the decoder to gather information across query-key pairs, effectively boosting performance. Further, TransMatcher is designed as a simplified decoder more suitable for image matching, which discards all attention implementations with softmax weighting and only keeps query-key similarity computation. This study demonstrates that Transformer can be effectively adapted to image matching and metric learning tasks with strong potential, and now it has been widely used in later research [163], [164] to improve the generalizability.

TABLE 2: Comparison of state-of-the-art supervised and unsupervised methods based on CNN and Transformer on two widely used datasets Market-1501 [4] and MSMT17 [24] in Person Re-ID. The performance of different pre-training conditions is reported. The supervised TransReID-SSL results are obtained by basic Transformer baseline [32] fine-tuning and the unsupervised TransReID-SSL results are obtained by Cluster-Contrast [112] fine-tuning. TransReID-SSL* refers to the results reproduced as a baseline in our experiments.

Method	Pre-training Conditions				Market-1501		MSMT17	
	Venue	Backbone	Data	Supervision	mAP	Rank-1	mAP	Rank-1
State-of-the-art methods for supervised Re-ID								
AGW [3]	TPAMI 2021	CNN	ImageNet	Supervised	87.8	95.1	49.3	68.3
CDNet [157]	CVPR 2021	CNN	ImageNet	Supervised	86.0	95.1	54.7	78.9
TransReID [32]	ICCV 2021	Transformer	ImageNet	Supervised	89.5	95.2	69.4	86.2
PHA [144]	CVPR 2023	Transformer	ImageNet	Supervised	90.2	96.1	68.9	86.1
TransReID-SSL [41]	arxiv 2021	Transformer	LUPerson	SSL	93.2	96.7	75.0	89.5
State-of-the-art methods for unsupervised Re-ID								
ICE [158]	ICCV 2021	CNN	ImageNet	Supervised	82.3	93.8	38.9	70.2
ISE [15]	CVPR 2022	CNN	ImageNet	Supervised	85.3	94.3	37.0	67.6
Cluster-Contrast [112]	ACCV 2022	CNN	ImageNet	Supervised	83.0	92.9	31.2	61.5
Cluster-Contrast [112]	ACCV 2022	CNN	LUPerson	SSL	84.0	94.3	31.4	58.8
PASS [159]	ECCV 2022	Transformer	LUPerson	SSL	88.5	94.9	41.0	67.0
TransReID-SSL [41]	arxiv 2021	Transformer	LUPerson	SSL	89.6	95.3	50.6	75.0
TransReID-SSL* [41]	arxiv 2021	Transformer	LUPerson	SSL	89.9	95.2	48.2	72.8
UntransReID (Ours)	-	Transformer	LUPerson	SSL	90.7	95.7	51.1	75.7

In addition, researchers try to investigate the generalization ability of Transformer in Re-ID. Ni *et al.* [163] employed different Transformers and CNNs as the backbone to assess the cross-domain performance from Market [4] to MSMT [24]. The results indicate that Vision Transformers outperform CNNs significantly. On this basis, a proxy task is introduced which mines local similarities shared by different IDs based on part aware attention, to promote the Transformer to learn generalized features without using ID annotations.

3.3 Transformer in Cross-modal Re-ID

In this subsection, we summarize the Transformer progress of three types of cross-modal problems that have received more attention in Re-ID: visible-image § 3.3.1, text-image § 3.3.2 and sketch-image § 3.3.3. Recently, Transformer has made a lot of novel works and influential breakthroughs in multi-modal learning in the field of vision. The key advantage is that the Transformer input consists of one or more token sequences and the attributes of each sequence, which allows the association of different modalities through attention mechanism without architectural modification [37].

3.3.1 Visible-Infrared Re-ID

visible-infrared Re-ID is a cross-modal retrieval task that aims at matching the daytime visible and nighttime infrared images [95], [165]. The major challenge of visible-infrared Re-ID is the modality gap between two types of images. The application of the Transformer provides many benefits to the visible-infrared Re-ID problem. For example, Transformers tend to learn shape and structure information, while CNNs rely on local texture information [39]. Due to the lack of colors and lighting conditions in infrared images, Vision Transformer can better capture modality-invariant

information and has stronger robustness. On the other hand, the vision transformer structure and attention enable local cross-modal associations to be easily established at the patch token level, which is essential for fine-grained properties of Re-ID, especially under large modality gap.

The mainstream approach in existing visible-infrared Re-ID is to learn modality-shared features, which decouple features into modality-specific and shared-modal features, and then focus on modality alignment at the feature level. Jiang *et al.* [166] start trying to adopt Transformer’s encoder-decoder architecture for modal feature enhancement and compensation to promote better alignment of RGB and IR modalities. They separately construct two sets of learnable prototypes for RGB and IR modalities to represent global modality information. In the Transformer decoder, the IR prototype is regarded as a query for the RGB modality and the part features at the token level of the RGB samples are used as keys and values. Compensated IR modal features are obtained by aggregating partial features through the correspondence of cross-attention between partial features and modal prototypes. In contrast, Liang *et al.* [167] introduce learnable embeddings to mine modality-specific features in Transformer in a manner similar to positional encoding, and employed a modality removal process to subtract the learned modality-specific features.

Considering the specificity of the body part position of the person object, Chen *et al.* [168] believe that position interaction can discover the underlying structural relationship between regions and provide more stable invariance for pose changes. A structure-aware position transformer (SPOT) is proposed to extract modality-shared representations. It exploits the attention mechanism to learn structure-related features guided by human key points and adaptively combines partially recognizable cues by modeling context

and position relations through a transformer encoder. Additionally, Feng *et al.* [169] focus on the interaction of local features across modalities, where they leverage attention to enrich the feature representation of each patch token by interacting with patch tokens from other modalities. Yang *et al.* [170] argue that each token in the self-attention mechanism in ViT is connected to a class token, where the attention score can be intuitively interpreted as a measure of token importance. To better align features of different modalities, they select top-k important visual patches from each attention head for localizing important image regions. Focusing on the modal invariant information of shallow features such as texture or contour information, Zhao *et al.* [171] utilize Transformer to encode the spatial information of each convolution stage of CNNs to fuse shallow and deep features to enhance the representation.

3.3.2 Text-Image Re-ID

Text-Image Re-ID refers to a cross-modal retrieval task, which aims at identifying the target object (person or vehicles) from an image gallery based on a given textual query, describing the target appearance [25], [172].

CLIP in Re-ID. As a milestone work of Transformer in multimodal applications, the proposal of Contrastive Language-Image Pre-training (CLIP) [173] opened up a new era of large-scale pre-training for text-image communication. CLIP uses text information to supervise the self-training of vision tasks, which turns a classification task into an image-text matching task. During the training process, a two-stream network including an image encoder and a text encoder processes text and image data respectively, and contrastive learning is used to learn the matching relationship between text-image pairs. The pre-trained model directly performs zero-shot image classification without any training data, achieving comparable supervised accuracy. Recently, CLIP has become a powerful tool for downstream text-image Re-ID tasks. Some Re-ID works directly introduce CLIP as a pre-training model with good generalization [174] or further expand the design of cross-modal association mining [40], [175], [176], and some works focus more on the utilization of Re-ID related textual information in CLIP to better assist downstream Re-ID tasks [177], [178].

Considering the effectiveness of directly fine-tuning CLIP, Yan *et al.* [175] explore the transfer of CLIP models to text-image Re-ID. Based on the pre-trained CLIP, they further capture the relationship between image patches and words to build fine-grained cross-modal associations. Inspired by CLIP, Zuo *et al.* [176] propose a language-image pre-training framework PLIP that is more suitable for person objects. To explicitly establish fine-grained cross-modal relations, a large-scale person dataset constructed with stylish generated text descriptions is proposed and three pretext tasks are introduced. The first is semantic-fused image coloring, which recovers the color information of gray-scale person images given a textual description. The second is visual-fused attributes prediction, which predicts masked attribute phrases in text descriptions through paired images. The last is visual-language matching. Instead, with CLIP as the initialization model, IRRA [40] designs a cross-modal implicit relation reasoning module to efficiently construct the relation between visual and textual represen-

tations through self-attention and cross-attention mechanisms. This fused representation is used to perform masked language modeling (MLM) task without any additional supervision and inference costs, achieving the purpose of effective inter-modal relation learning. Compared with the previous method which utilizes single-modal pre-trained external knowledge and lacks multi-modal corresponding information, these CLIP-based text-image Re-ID methods have achieved significant performance improvements.

Compared with the one-hot label of image classification, CLIP-Re-ID [177] demonstrates that more detailed image text descriptions can help the visual encoder learn better image features, especially for fine-grained tasks such as Re-ID that lack precise descriptions. Inspired by the learnable prompt used by CoOp [179], CLIP-Re-ID designs a two-stage training strategy. It combines ID-specific learnable tokens to give ambiguous textual descriptions in the first stage and these tokens together with the text encoder provide constraints for optimizing the image encoder in the second stage. Besides, Wang *et al.* [178] focused on exploring the potential of text in the vision-language pre-training model on the text-image Re-ID baseline based on CLIP. They proposed two strategies for describing the corpus to be integrated into the baseline: one is to rearrange fine-grained information in attribute phrases to generate correct but different types of damaged or additional sentences, and the other constructs unified attribute descriptions from core attributive noun-pair class-related prompt templates.

3.3.3 Sketch/Skeleton Re-ID

Sketch-to-photo Re-ID represents a cross-modal matching problem whose query sets are sketch images provided by artists or amateurs [113], while the query images in skeleton Re-ID are generated by pose estimation [42]. These two tasks share similar spirits in large information asymmetry.

Sketch-image Re-ID. The significant difference in region-level information between sketches and images is a challenge due to the abstraction and iconography of sketch images. The correlation between object shape and local information plays an important role for sketch-image Re-ID. The advantage of Transformer in learning global-level feature representations shows excellent discriminative ability in sketch photo recognition [119], [141]. Chen *et al.* [141] experimentally verified that the method using ViT as the backbone has a significant performance improvement over most CNN-based sketch-image Re-ID methods. Therefore, they construct a strong baseline based on vision transformer for sketch-image Re-ID. In order to narrow the gap between sketches and images, Zhang *et al.* [180] designed a token-level cross-modal exchange strategy in Transformer under the guidance of identity consistency to learn modality compatible features. Local tokens of different modalities are classified into different groups and assigned specific semantic information to construct a semantically consistent global representation.

In particular, a new, challenging and modally agnostic person re-identification problem has recently been proposed [49]. It comprehensively considers descriptive queries such as supplementary text or sketch modalities for general images to achieve multi-modal unified re-identification. Benefiting from CLIP, UNIRe-ID [49] developed a simple dual-

encoder transformer architecture for multimodality feature learning and designed a task-aware dynamic training strategy that adaptively adjusts the training focus according to the difficulty of the task. This work demonstrates the power of Transformer in multi-modal learning and also opens up the direction for the future promotion of Re-ID.

Re-ID with Skeleton Data. Person Re-ID via 3D skeletons differs from traditional Re-ID, which relies on visual appearance features such as color, outline, etc., and mainly utilizes the 3D locations of key body joints to model unique body and motion representations. TranSG [42] is proposed as a general Transformer paradigm for learning feature representations from skeleton graphs for person Re-ID. Its core idea is to model the 3D skeleton as a graph and use Transformer for full-relational learning of body joint nodes, which simultaneously aggregates key relationship features of body structure and motion into a graph representation.

3.4 Transformer in Special Re-ID Scenarios

We investigate that the vision Transformer is also applied to some more open and complex task settings in Re-ID. The more special Re-ID types than the scenarios mentioned above are summarized in this subsection for discussion.

Cloth-changing Re-ID. It is a challenging Re-ID task in long-term scenarios where the person may change the cloths in an unknown pattern. Cloth-changing Re-ID is a challenge unique to persons, which is a difficult but more practical problem [116]. In this scenario, the discriminative feature representation dominated by the visual appearance of clothing will be invalid. Existing research in tackling this more intricate challenge has initiated initial investigations into the application of Transformers. Lee *et al.* [181] evaluate different backbones in the cloth-changing Re-ID scenario, and Transformer demonstrated notable performance advantages when compared to CNNs. On this basis, in order to further eliminate the influence of characteristics related to clothing or accessories, an attribute de-biasing module is designed. The core idea is to use the generated attribute labels for person instances as auxiliary information and adopt a gradient reversal mechanism based on adversarial learning to learn attribute-agnostic representations.

Human-centric Tasks. The success of the general large model built by Transformer lies in its ability to handle multiple tasks. Recent work has attempted to focus on human-centric general model design to facilitate the Re-ID task. Human-centric perception integrates visual tasks such as pedestrian detection, pose estimation, attribute recognition, and human parsing. Person Re-ID is one of the human-centric tasks. These tasks all have in common that they rely on the basic structure of the human body and the properties of body parts. Previous studies have experimentally verified that training human-centric tasks together can benefit each other [182]. It is challenging to unify large-scale multiple tasks into a scalable model due to the different structure and granularity of annotations and expected outputs of different tasks requiring separate output headers for each task. UniHCP [182] presents a flexible Transformer encoder-decoder structure to avoid task-specific output heads. The core idea is to define task-specific queries in the decoder and design a task-guided interpreter to interpret each query

token independently. Outputs of the same modality share the same output unit, enabling maximum parameter sharing among all tasks while learning human-centric knowledge at different granularities. Additionally, Tang *et al.* [183] establish a large-scale dataset HumanBench, specifically designed for human-centric pre-training. To address task conflicts arising from diverse annotations in supervised pre-training, a projector assisted hierarchical pre-training method is proposed. The core idea involves constructing a hierarchical structure: sharing the weights of the backbone across all datasets, while restricting the weights of the projector to be shared only among datasets of the same tasks, and the weights of the head to be shared for a single dataset.

On the other hand, unlabeled person data is plentiful, and researchers use Transformer’s strong scalability for large-scale self-supervised training to learn human-centric representations. Considering methods such as contrastive learning or masked image modeling that failed to explicitly learn semantic information, SOLIDER [184] uses Transformer to generate pseudo semantic labels for every token based on prior knowledge of human images and introduces a token-level semantic classification pretext task to learn a stronger human semantic representation. With the mutual promotion of large-scale learning in multiple human-centric tasks, SOLIDER can achieve superior performance compared with the state-of-the-art unsupervised pre-training methods [41], [159] for Re-ID, which can be regarded as a further advance in the Re-ID community.

Person Search. Person search is an end-to-end method that aims to jointly solve the two sub-problems of detection and person Re-ID using more efficient multi-task learning methods. Since the goals of person detection and person Re-ID are conflicting and it is difficult to jointly learn a unified feature representation, Yu *et al.* [185] propose to decompose feature learning into successive steps in the T stage of a multi-scale Transformer to gradually learn from coarse to fine embedding. Unlike some existing multi-scale Transformers that learn different scale information based on patches of different sizes, they leverage a series of convolutional layers with different kernels to generate multi-scale tokens. Furthermore, to produce more occlusion-robust representations, they design to exchange partial tokens of instances in mini-batches, and then compute occlusion attention based on mixed tokens. Different from the shuffling and regrouping strategy in TransReID [32], they for partial tokens in a single instance. PSTR [33] is also designed as a person search multi-scale learning scheme of the Transformer architecture. It develops a PSS module consisting of a detection encoder-decoder and a discriminative re-identification decoder, where the detection encoder-decoder employs backbone features and three cascaded decoders are employed. The Re-ID decoder takes a feature query from one of the three detection decoders as input, and a multi-level supervision scheme is designed to provide different input Re-ID feature queries and box sampling locations. In order to achieve multi-scale expansion, the features of different layers use PSS modules and are concatenated to perform instance-level matching with queries.

Group Re-ID. It is a group-level Re-ID by using the contextual information to match a small number of individuals in a group together [186]. Since people usually have

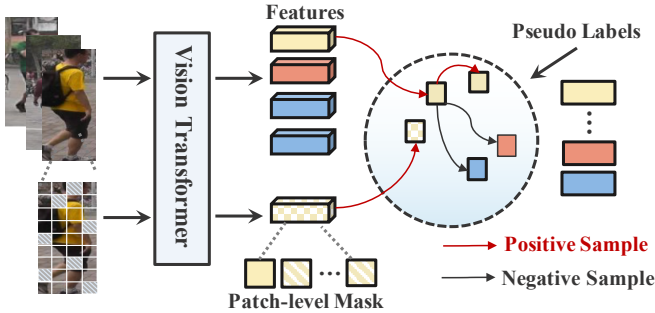


Fig. 5: The proposed unsupervised Transformer baseline for Re-ID enhanced with a patch-level mask learning strategy.

group and social attributes, group actions are preferred in most real-world scenarios. Group Re-ID has gradually attracted the attention of researchers, which needs to deal with challenges such as membership and layout changes. Existing group Re-ID methods are mainly based on the combined framework of CNN and GNN. However, these structures are deficient in position modeling and have weak ability to describe group layout characteristics. Inspired by the position embedding in the transformer, Zhang *et al.* [187] design the second-order Transformer model SOT to deal with the layout features in group Re-ID. It consists of intra-member and inter-member modules, where each member in the group image is first cropped, and then each member is segmented into multiple sub-patches. The intra-member module extracts first-order labels as per-member features by modeling the relationship between sub-patches through a transformer. The member-to-member module models the relationship between members through uncertainty and extracts second-order tokens through transformers.

Re-ID in UAVs. Object Re-ID in UAVs involves identifying specific objects within a multitude of aerial images captured from a dynamic bird’s-eye view [188]. It also has broad application prospects in different scenarios. Re-ID using aerial images captured by UAVs is an under-explored scenario. Unlike the widely used fixed cameras, the images captured by UAVs are more complex than fixed city cameras. Unavoidable continuous rapid movement and height changes lead to large differences in image viewing angles. Chen *et al.* [21] analyzed the vehicles and pedestrians in the bird’s-eye view and concluded that Re-ID faces two key challenges in UAV scenarios: bounding boxes with significant size differences and objects with uncertain rotation direction. They designed a novel idea of a feature-level rotation strategy by taking advantage of the Transformer structure. Benefiting from the corresponding relationship between image patches and token-level features in Vision Transformer, the insight of this work is to simulate rotation operations on the initially learned patch features to generate enhanced diversity rotation features. Compared with image rotation data augmentation, this design alleviates the loss of key information caused by cropping or padding when the image is rotated. Another study [189] explores the enhancement of the Pyramid Vision Transformer (PVT), leveraging multi-scale features for object Re-ID in UAV scenarios.

TABLE 3: Evaluation results of our Transformer-based visible-infrared cross-modal unsupervised Re-ID baseline on two datasets RegDB [191] and SYSU-MM01 [95].

Method	RegDB				SYSU-MM01			
	V-T		T-V		All Search		Indoor Search	
	mAP	R1	mAP	R1	mAP	R1	mAP	R1
OTLA ECCV22 [192]	29.7	32.9	28.6	32.1	27.1	29.9	38.8	29.8
ADCA MM22 [50]	64.1	67.2	63.8	68.5	42.7	45.5	59.1	50.6
ACCL CVPR23 [18]	65.4	69.5	65.2	69.9	51.8	57.3	62.7	56.2
UntransReID (Ours)	69.9	76.3	69.3	76.8	52.5	51.9	66.0	57.5

4 NEW UNSUPERVISED TRANSFORMER BASELINE

After conducting a thorough review of Transformer’s work in Re-ID, we are confident that large-scale pre-trained Transformers hold substantial promise for unsupervised Re-ID and warrant further exploration. Most previous Re-ID works are pre-trained on ImageNet, due to the lack of large-scale person datasets. In fact, pre-training on person datasets, such as LUPerson [58], is better suited for Re-ID task and aligns with future development trends. Our survey reveals that some studies [41], [159] verify the evident advantages of using Transformer pre-training on LUPerson. In order to further promote the progress of the Re-ID community, we propose a single/multi-modal general unsupervised Re-ID baseline. Specifically, our baseline follows the Re-ID method [112] of contrastive learning of pseudo-labels generated by clustering, and uses the TransReID-SSL [41] pre-trained Transformer as a powerful initialization model. On this basis, leveraging the characteristics of Transformer, we have devised the following design for UntransReID.

Single-modal Unsupervised Re-ID. Inspired by existing Transformer-based masked image modeling self-supervised methods [122], [190], we design a patch-level mask enhancement strategy integrated into the unsupervised training process. Our core idea is to adopt a series of learnable tokens to mask part of the image patches as an augmentation and establish the relationship between the original features and the mask features during the training process as a supervisory signal to guide model learning. On the other hand, aligning the mask features with the original features inherently encourages the model to learn local fine-grained information. For input images, we define the set $\mathcal{X}^g = \{x_i^g | i = 1, 2, \dots, n\}$ and set $\mathcal{X}^l = \{x_i^l | i = 1, 2, \dots, n\}$ respectively as the original input and mask enhancement input. Patch embedding operations are used to get preliminary tokens $x_i^g \in \mathbb{R}^{N \times D}$ and $x_i^l \in \mathbb{R}^{N \times D}$, where N and D represent the number of patches and the dimension of the token. We initialize a set of learnable mask tokens $\mathcal{M}^l = \{m_i^l | i = 1, 2, \dots, m\}$, randomly replacing p of the tokens in \mathcal{X}^l as the final input. The corresponding output class tokens after Transformer model learning are $\{f_i^g | i = 1, 2, \dots, n\}$ and $\{f_i^l | i = 1, 2, \dots, n\}$. We calculate the contrastive loss between the original image features and mask-enhanced features as:

$$\mathcal{L} = -\log \frac{\exp(f_i^g \cdot f_i^l / \tau)}{\sum_{j=1}^k \exp(f_i^g \cdot f_j^l / \tau)}, \quad (2)$$

where k represents the batch size and f_j represents the original features within the batch.

Cross-modal Unsupervised Re-ID. For the transformer-based unsupervised visible-infrared cross-modal Re-ID, we devise a dual-path transformer that adopts two modality-specific patch embedding layers and a modality-shared transformer. Each modality-specific patch embedding layer comprises an IBN-based Convolution Stem (ICS) [41] to capture modality-specific information. The modality-shared transformer is introduced to learn a multi-modality sharable space. On the basis of [112], two modality-specific memories are constructed for mining inter- and intra-class information within each modality with contrastive learning. To further ensure the modality generalization capability, we adopt random channel augmentation following [17] as an extra input to the visible stream for joint learning.

Analysis of Results. Tab. 2 and Tab. 3 respectively present the evaluation results for our single-modal and cross-modal unsupervised Re-ID baselines. For single-modal unsupervised Re-ID, the structural characteristics of the Transformer allow us to generate augmented samples by applying local masks at the patch level, enabling the construction of supervisory signals. On the powerful Transformer backbone pretrained on LUPerson [41], our baseline combined with the enhancement strategy and contrastive learning [112], achieves performance comparable to state-of-the-art methods. For cross-modal Re-ID, existing state-of-the-art methods are all based on CNNs and require complex cross-modal association designs, whereas our Transformer baseline achieves state-of-the-art performance with a simple design across several infrared-visible Re-ID datasets.

5 ANIMAL RE-IDENTIFICATION

In addition to objects such as persons and vehicles, which are currently widely studied in the field of Re-ID, the demand for wildlife protection makes animal re-identification gradually attract the attention of researchers [196], [197]. Animal Re-ID has important applications in many fields such as ecology, conservation biology, environmental monitoring, popular science education, and agriculture. It provides a powerful tool for understanding nature, protecting biodiversity, and maintaining ecological balance. Considering that the existing Re-ID surveys are all for persons or vehicles, our survey will cover a wider range of Re-ID objects to promote Re-ID research. Animals are an important part of Re-ID objects. This section summarizes related papers on animal Re-ID in recent years.

Different from the mature development of person and vehicle re-identification technology, animal re-identification technology is still in a relatively early stage. One of the most intuitive challenges is that the uncontrollable environmental factors in the wild make it very complicated to collect animal images and label their individual information. As shown in Fig. 6, compared to humans, animal species exhibit a diverse array of unique characteristics across different species. The core problem of animal individual re-identification is to mine and analyze the discriminative features specific to a species. Our survey summarizes the animal Re-ID datasets released in recent years in § 5.1. In addition, deep learning-based animal re-identification techniques are introduced in § 5.2.

5.1 Animal Re-ID Datasets

Due to the diversity of environments and ways in which different animals live, the collection of animal data is not as simple as that of persons and vehicles. In addition to using surveillance cameras and ordinary digital cameras, camera traps, drones, and infrared thermography are also important devices. Therefore, most animal Re-ID datasets focus on annotating identity information in order to complete specific individual recognition, without a clear definition of camera. The emergence of more and more animal Re-ID datasets in recent years has promoted the research progress [9], [198], [199]. As shown in Tab. 4, we provide an overview of animal Re-ID datasets in recent years. We present key features of different species in Re-ID, which are also special challenges in the animal Re-ID task. In addition, some animal data such as long-lived species are collected over a period of several years. This kind of data with a long time span provides more comprehensive information and supports deeper analysis for long-term ecological research and species protection [200]. This also implies that Re-ID will be more challenging since the appearance of animals can change dramatically over time. We also report the time span for each dataset.

Our survey collects animal Re-ID datasets from different sources in recent years to promote Re-ID research, some of which are paper publications [193], [194], [198], [199], [210], and some in the form of competitions [9], [195]. This is mainly due to the fact that the datasets of many papers are not publicly available, and datasets from diverse sources allow advanced research to be conducted on them. First, some domestic or laboratory animal datasets are introduced. Bergamini *et al.* [201] collect cattle head images in farms for Re-ID, considering that the cattle heads can show sufficient texture, shape and patch characteristics. Li *et al.* [204] shot at a real cattle farm and built a dataset of 13 cows. The Cows2021 dataset [205] contains 186 Holstein-Friesian cattle, which took a month to capture from a bird’s-eye view on the farm. For cattle, their personalized black and white coat pattern patches are an important distinguishing characteristic of individuals. YakRe-ID-103 [198] is a Yak dataset of highland pasture scenes. Yaks mostly have black fur and are typically texture-less animals, making it difficult to distinguish individuals. The most unique features are the thickness, bending and direction of the horn. Besides, a dataset of six zebrafish recorded in a laboratory setting is presented by [203]. They propose reliable re-identification through the stripes of zebrafish from the side view.

Re-ID for wild animals is relatively more challenging due to complex animal habits and uncontrollable environments. ATRW [9] is a dataset containing 92 Amur tigers collected in multiple large wild zoos with bounding boxes, pose key points and identity annotations. Korschens *et al.* [193] collected images of forest elephants in national parks and constructed a data set containing 2078 images of 276 elephant individuals. The dataset spans approximately 15 years, which reflects some of the aging effects and dramatic changes in the physique of elephants. The most identifiable tusks and scars of elephants also change over time, exacerbating the difficulty of Re-ID. Chan *et al.* [207] constructed a short-term and long-term dataset for honeybee re-identification. They mainly focused on the abdomen of

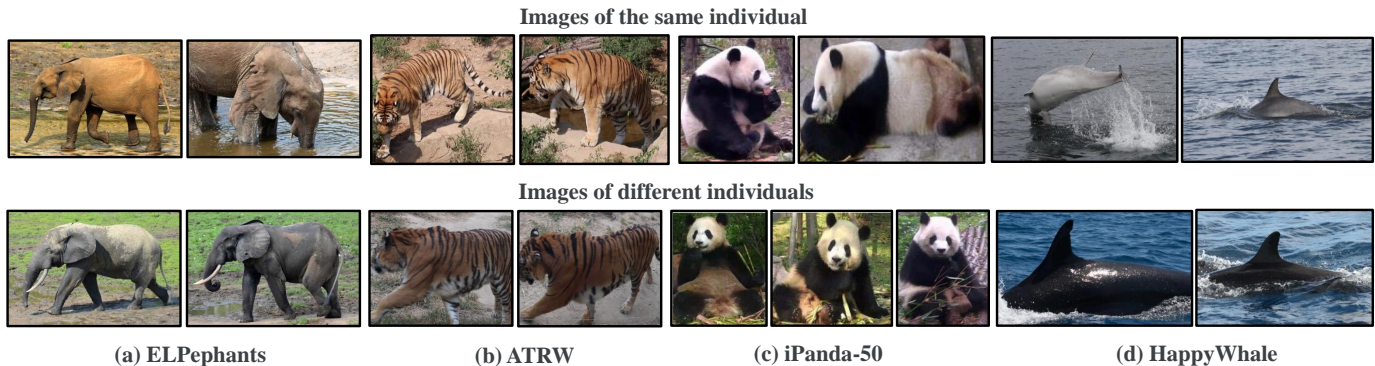


Fig. 6: **Images of different species in Animal Re-ID.** Unlike the widely studied person and vehicle Re-ID, animal individuals of the same species have extremely similar appearances. Different species have their own unique discriminative characteristics, such as (a) the tusks and injury marks of elephants [193], (b) the coat pattern of amur tiger [9], (c) the eyes of giant pandas [194], (d) the dorsal fin, back, and flank of whales [195].

TABLE 4: **Summary of animal Re-ID datasets from recent years.**

Dataset	Species	IDs	Images	Time	Feature	Span	Source	Available
CattleRe-ID [201]	cattle	-	-	2018	face	-	farm	×
DolphinRe-ID [202]	dolphin	185	3544	2018	fin	12 years	-	×
Elpephants [193]	elephant	276	2078	2019	body, tusk	15 years	national park	×
ATRW [9]	Amur tiger	92	3649	2019	stripe	-	wild zoos	✓
zebrafishRe-ID [203]	zebrafish	6	2224	2020	side view	-	laboratory	✓
CowRe-ID [204]	cow	13	3772	2021	coat pattern	-	farm	✓
YakRe-ID-103 [198]	yak	103	2247	2021	horn	-	highland pastures	×
Cows2021 [205]	cattle	182	13784	2021	coat pattern	1 month	farm	✓
iPanda-50 [194]	giant panda	50	6874	2021	local	-	Panda Channel	✓
SealID [199]	seal	57	2080	2022	pelage pattern	10 years	Lake Saimaa	✓
FiveVideos [196]	pigeon, fish, pig	93	20490	2022	-	-	Pixabay	✓
BelugaID [206]	beluga whale	788	5902	2022	scarring pattern	4 years	Cook Inlet	✓
Honeybee [207]	honeybee	181	8962	2022	abdomen	multiple weeks	colony entrance	×
HappyWhale [195]	30 species	15587	51033	2022	fin, head, flank	very long	28 organizations	✓
SeaTurtleID [200]	sea turtle	400	7774	2022	-	12 years	Laganas Bay	✓
LeopardID [208]	African leopard	430	6795	2022	spot pattern	11 years	-	✓
HyenaID [209]	spotted hyena	256	3104	2022	spot pattern	-	-	✓
PolarBearVidID [210]	polar bear	13	138,363	2023	-	-	zoo	✓
Wildlife-71 [211]	71 species	≈2059	≈108,808	2023	-	-	internet	✓

honeybee for individual discrimination, and the time span of the long-term dataset reached 13 days. iPanda-50 [194] is a giant panda Re-ID dataset collected through giant panda streaming videos, which contains 50 giant pandas of different ages including cubs, juveniles, and adults. The Saimaa ringed seal is an endangered subspecies found only in Lake Saimaa, Finland. Individual ringed seals have unique fur patterns, and individual re-identification is of great value in monitoring endangered animals [199]. SealID [199] is a benchmark for Saimaa ringed seal re-identification, which takes into account challenges such as the deformable nature of seals and low contrast between the ring pattern. SeaTurtleID [200] is a large-scale dataset containing images of sea turtles captured in the wild. This dataset is time-stamped and spans up to 12 years. Considering the impact of timestamps for unbiased evaluation of animal Re-ID methods, the dataset also provides time-aware partitioning

of reference and query sets. Happywhale [195] is an open source platform to facilitate the identification of individual marine mammals. PolarBearVidID [210] provides a video-based dataset of polar bears, which is challenging that individuals lack significant unique visual features. While the majority of existing datasets are tailored to a single species, recent research has introduced large-scale Re-ID datasets that encompass multiple species. The Wildlife-71 dataset [211], proposed as a dataset that aggregates existing datasets and partial web data, includes Re-ID data from 71 different wildlife categories. Indeed, it is observed that many existing animal datasets are relatively small in scale, and aggregating multi-species data proves advantageous for deep learning technologies. In addition, we observe that the animal Re-ID is much less explored compared with other objects. Recognizing the animals suffer from additional severe occlusions and viewpoint changes compared with persons.

5.2 Animal Re-ID Methods

In this survey, we mainly focus on advanced deep learning methods for animal Re-ID due to their powerful performance compared with other traditional solutions. Our survey broadly categorizes these methods into three groups: learning with global animal images, learning with key local body areas, and learning with auxiliary information. Note that Transformer is seldom explored in this area.

Global Image Based Methods. Many existing studies draw upon the conventional approaches of person Re-ID, directly feeding entire animal images into deep neural networks to acquire reliable feature representations [202]. Taking cues from person Re-ID methodologies like local maximal occurrence [212], Bruslund *et al.* [203] introduce two feature descriptors consisting of color and texture to reliably re-identify zebrafish from side-view. Considering that the patterns on manta rays are usually in uncertain positions, Moskvayak *et al.* [213] devise a loss function to minimize the distance between the same individual observed from various viewpoints, guiding the learning of pose-invariant features. Porrello *et al.* [214] propose a general view knowledge distillation method for Re-ID tasks. The core idea is to use the diversity of the target in different views as a teaching signal, allowing students to use fewer views to restore it and learn more robust features. For giant panda re-identification, Wang *et al.* [194] design a multi-stream structure to learn local and global features. In order to mine local fine-grained information from the global image, a patch detector is adopted to automatically capture the most discriminative local patches without additional part annotations.

Local Area Based Methods. Among the related work on animal Re-ID, some research focuses on specific parts of the animal. They extract the most discriminative areas of the original image during the data collection stage, such as the head of a cow [201], elephant ears [215], whale tails [216], dolphin fins [202], [217], [218], etc. Bergamini *et al.* [201] employ CNNs for direct feature extraction from self-collected cattle head datasets and utilize KNN (k-nearest neighbors) for classification. For fine-grained images of elephant ears and whale tails, Weideman *et al.* [215] design their approach to extract boundary information from color and texture transitions, along with intensity variations, to effectively discern the outlines of critical regions.

Auxiliary Information Based Methods. Zhang *et al.* [198] utilize a simplified definition of the pose of the yak’s right or left head as an auxiliary supervision signal to enhance feature learning. Li *et al.* [9] employed the results of pose key point estimation to model the tiger image into 7 parts including the trunk, front legs, and hind legs to learn local features. In order to learn the unique body markings of animal individuals with similar appearance, Moskvayak *et al.* [219] proposed a heat map enhancement method to display the location information of introduced animal landmarks in the Re-ID model. When dealing with species exhibiting similar pelage or fur patterns, Nepovinskykh *et al.* [220] employed the Sato tubeness filter to extract the fur pattern from the image, mitigating the impact of interfering factors like lighting. Siamese networks [221] trained with triplet loss are used for subsequent matching.

5.3 A Unified Benchmark for Animal Re-ID

In fact, existing deep learning-based animal Re-ID methods are still in the early stages of development, and we generally summarize their main limitations: (1) *Unclear Task Boundaries.* Many animal-related studies do not have clear task definitions, some of which are regarded as fine-grained recognition or individual classification [194]. They typically concentrate solely on distinguishing between different individuals and pay little attention to whether they can be reliably re-identified across various settings. However, in this survey, we emphasize animal re-identification, with the goal of accurately identifying the same individual across different timeframes, environments, or viewpoints. (2) *Limited method applicability.* Many existing methods leverage the distinctive traits of particular species to develop their approaches, with some specifically focusing on curating datasets for certain body parts of the animals [201], [207], [217], [220]. These approaches prove challenging to adapt for broader application in Re-ID across different species and exhibit limited scalability. (3) *Inconsistent experimental settings.* Existing animal Re-ID methods adopt varying experimental settings. Some of the work is conducted experimentally in a closed-world setting, which involves identifying objects within known and limited categories. In most cases, the Re-ID system is not aware of all possible categories during training, necessitating its ability to handle unforeseen categories. Some research has also been conducted in scenarios that better align with the open-set nature of Re-ID tasks. This makes performance comparison between different methods a challenging task.

To further advance research and realize the full potential of animal re-identification for practical applications, it is critical to establish standardized benchmarks and develop more robust, scalable techniques. In this survey, we conducted extensive animal Re-ID experiments using multiple state-of-the-art general Re-ID methods to address the aforementioned issues. Our work in this section covers a unified evaluation setting, a comparison of different backbone methods, and an analysis of the Transformer’s suitability for animal Re-ID. The code will be publicly available.

5.3.1 Animal Re-ID experiments.

Datasets. We chose datasets featuring various species, such as giant pandas [194], elephants [193], seals [199], giraffes [222], zebras [222], leopards [208] and tigers [9] for our evaluation. Since the datasets only provide original images and corresponding identity annotations, we uniformly divide them into training sets and test sets for the Re-ID task. Specifically, we divide each data set into 70% of all identities as training data, and the remaining 30% as test data. Ensure that the identities of the test set have not appeared in the training set. In the testing phase, we regard each image in the test set as a query, and all images in the test set except the query image constitute the gallery. The results for FiveVideos in Tab. 5 are obtained using only pig data. The results for GZGC-G and GZGC-Z are using giraffe data and zebra data, respectively.

Evaluation Metrics. The performance is evaluated by two widely used metrics in Re-ID tasks: Cumulative Matching Characteristic (CMC) and the mean Average Precision

TABLE 5: Evaluation results of Re-ID methods on multiple animal datasets.

Dataset	BoT [69]			TransReID [32]			RotTrans [21]		
	mAP	R1	mINP	mAP	R1	mINP	mAP	R1	mINP
iPanda-50	28.4	72.5	9.8	37.9	88.8	10.5	42.6	91.7	12.9
ELPephants	15.8	32.3	4.5	15.3	40.2	3.1	30.2	56.0	9.7
SealID	49.1	82.2	7.2	42.6	82.8	6.3	48.3	83.5	7.4
ATRW	65.2	98.4	32.5	64.1	98.3	33.0	66.9	97.9	35.4
GZGC-G	47.4	46.7	38.4	49.1	48.9	39.0	48.9	47.8	40.4
GZGC-Z	13.7	23.5	5.6	16.3	26.0	7.4	16.2	26.7	7.2
LeopardID	27.3	60.1	9.9	31.6	63.7	12.5	32.5	63.0	13.3

(mAP). It’s worth noting that in the context of person and vehicle Re-ID, correctly matched objects captured by the same camera are typically excluded from the evaluation, while only objects captured by different cameras are considered. However, in animal Re-ID, where explicit camera information is often lacking in most animal datasets, we calculate all correctly matched objects uniformly. In many cases, simple samples with small viewing angle changes will lead to high Rank-k accuracy. Therefore, we calculate a new metric mean Inverse Negative Penalty (mINP) [3], which reflects the cost of finding the hardest matching sample.

Analysis of Results. To evaluate the performance of different backbones in animal Re-ID, two Re-ID methods that are generally applicable to various objects, CNN-based BoT [69] and Transformer-based TransReID [32], are employed in our experiments. As shown in Tab. 5, the mean accuracy of existing state-of-the-art Re-ID methods applied directly to animals is generally low. This also underscores that Animal Re-ID, distinct from the widely studied Re-ID objects, poses unique challenges and requires more targeted solutions in the future. The Transformer method performs better in most cases. In addition, considering some characteristics of animal Re-ID that are different from conventional person Re-ID such as camera view and diverse orientations, we choose a state-of-the-art Transformer-based method of object Re-ID in UAVs which is mentioned in § 3.4, RotTrans [21], for evaluation. We believe that images of different species (e.g., marine and terrestrial animals) will exhibit a variety of rotation angles rather than being in a standing position as persons. Consequently, RotTrans demonstrates superiority in most animal Re-ID scenarios as a method that helps to learn rotation-invariant representations. Recently, researchers have proposed the development of a Re-ID model capable of handling any unseen wildlife category [211]. The concept is similar to the domain generalization problem in conventional Re-ID tasks, and their solution involves leveraging larger scale and more diverse data. Differently, our benchmark is designed for the general animal Re-ID task, specifically aiming at the methodological level of multi-species applicability.

6 CONCLUSION AND FUTURE PROSPECTS

6.1 Under-Investigated Future Prospects

We discuss two under-investigated open-issues, including *Unified Large-scale Foundation Model for Re-ID* and *Efficient Transformer Deployment for Re-ID*.

Unified Large-scale Foundation Model for Re-ID. To meet the practical application demands of Re-ID, it primarily involves the utilization of a unified large-scale model that accommodates multi-modality and multi-object scenarios. In cross-modal Re-ID, existing research often concentrates on just two specific modalities. However, the sources of query cues are highly diverse, and exploring how to integrate information from various senses or data sources to create a genuinely modality-agnostic universal Re-ID model is a matter of significance. Transformer shows great potential in this problem, owing to its flexible handling of multi-modal inputs, robust relationship modeling capabilities, and scalability for processing large-scale data. The latest research reveals that Transformer continues to make breakthroughs in constructing multi-modal large models. For instance, Meta-Transformer [223] is proposed to understand 12 modal information and offer a borderless multi-modal fusion paradigm. Besides, unifying various tasks related to Re-ID that target the same objective is a development direction with practical significance. In our survey, it is evident that several recent studies have made breakthroughs by utilizing Transformer models to unify diverse vision tasks centered around humans [182], [183], [184]. Furthermore, constructing a universal model for multiple objects in Re-ID poses a significant challenge. This implies that many methods rely on specific information face difficulties. Particularly in animal Re-ID, the creation of a multi-species robust model holds great importance for practical applications.

Efficient Transformer Deployment for Re-ID. Our survey demonstrates that Transformers indeed exhibit powerful performance in the Re-ID field. However, due to the substantial computational support required for self-attention calculations, the associated resource consumption is relatively high. In practical applications, such as video surveillance and intelligent security, there is an increasing demand for real-time performance and lightweight deployment of Re-ID models [135]. Balancing the preservation of the Transformer’s robust performance in Re-ID with the imperative to reduce computational complexity becomes a crucial direction for future research. Besides, many pre-trained large-scale general foundation have been developed in general area, how to efficiently transfer the general knowledge to the specific Re-ID tasks is also worth studying [224]. considering the catastrophic forgetting problem in large-scale dynamic updated camera network, how to efficiently fine-tune the previously learned Re-ID models to downstream scenarios is another important direction to explore [225], [226].

6.2 Summary

From our survey, it is evident that in the past three years, Transformer has experienced rapid development in the Re-ID field, particularly demonstrating strong advantages in more challenging scenarios such as multi-modal and unsupervised settings. We provide an in-depth analysis of the advantages of Vision Transformer in four aspects, corresponding to four Re-ID scenarios:

- 1) *Transformer in Image/Video Based Re-ID*: At the backbone level, the Transformer entirely relies on the attention mechanism, providing it with universal modeling capabilities for global, local, and spatio-temporal relationships. This inherent capability facilitates the effortless

extraction of global, fine-grained, and spatio-temporal information in regular image and video Re-ID tasks.

- 2) *Transformer in Re-ID with Limited Data/Annotations*: The emergence of Transformer provides more possibilities for unsupervised learning. Beyond conventional discriminative learning approaches, such as contrastive learning, a broader spectrum of self-supervised paradigms (e. g. masked image modeling), has gained widespread attention and exploration. Furthermore, the Transformer exhibits superior adaptability to large-scale data, facilitating extensive self-supervised pre-training of more powerful and generalized models for addressing Re-ID with limited data or annotations.
- 3) *Transformer in Cross-modal Re-ID*: Transformer provides a unified architecture to effectively handle data of different modalities, especially the connection of vision and language. The multi-head attention mechanism possesses the capability to aggregate features across various feature spaces and global contexts, and the highly adaptable encoder-decoder structure is capable of accommodating diverse types of inputs and outputs. Consequently, the Transformer is particularly well-suited for establishing inter-modal associations and facilitating the fusion of multi-modal information in cross-modal Re-ID tasks.
- 4) *Transformer in Special Re-ID Scenarios*: Driven by the demands of practical applications, the Re-ID field has given rise to a range of specialized and challenging scenarios, such as cloth-changing Re-ID, end-to-end Re-ID, group Re-ID, Re-ID in UAVs, and human-centric tasks. The initial exploratory efforts of Transformer in tackling these intricate challenges have showcased remarkable scalability and adaptability.

This survey predominantly encompasses transformer-based Re-ID papers, primarily focusing on widely studied objects like persons and vehicles. Considering the application of Transformer in single-modal/cross-modal unsupervised Re-ID, which has not been fully explored by existing research, we present a Transformer-based baseline that achieves state-of-the-art performance on multiple single-modal/cross-modal Re-ID datasets. In particular, we explore the field of animal Re-ID, an area that continues to encounter challenges and unresolved issues. We develop unified experimental standards for animal Re-ID and evaluate the feasibility of employing Transformer in this context, laying a solid foundation for future research. Additionally, we delve into the future prospects of Transformers in Re-ID, aiming to further stimulate subsequent research.

REFERENCES

- [1] D. Gray, S. Brennan, and H. Tao, "Evaluating appearance models for recognition, reacquisition, and tracking," in *PETS*, vol. 3, no. 5, 2007, pp. 1-7.
- [2] C. C. Sun, G. S. Arr, R. P. Ramachandran, and S. G. Ritchie, "Vehicle reidentification using multidetector fusion," *IEEE TITS*, vol. 5, no. 3, pp. 155-164, 2004.
- [3] M. Ye, J. Shen, G. Lin, T. Xiang, L. Shao, and S. C. H. Hoi, "Deep learning for person re-identification: A survey and outlook," *IEEE TPAMI*, pp. 1-1, 2021.
- [4] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, "Scalable person re-identification: A benchmark," in *ICCV*, 2015, pp. 1116-1124.
- [5] E. Ahmed, M. Jones, and T. K. Marks, "An improved deep learning architecture for person re-identification," in *CVPR*, 2015, pp. 3908-3916.
- [6] A. Hermans, L. Beyer, and B. Leibe, "In defense of the triplet loss for person re-identification," *arXiv preprint arXiv:1703.07737*, 2017.
- [7] Z. Zhong, L. Zheng, D. Cao, and S. Li, "Re-ranking person re-identification with k-reciprocal encoding," in *CVPR*, 2017, pp. 1318-1327.
- [8] L. Zheng, Y. Yang, and A. G. Hauptmann, "Person re-identification: Past, present and future," *arXiv preprint arXiv:1610.02984*, 2016.
- [9] S. Li, J. Li, H. Tang, R. Qian, and W. Lin, "Atrw: a benchmark for amur tiger re-identification in the wild," *arXiv preprint arXiv:1906.05586*, 2019.
- [10] G. Wang, Y. Yuan, X. Chen, J. Li, and X. Zhou, "Learning discriminative features with multiple granularities for person re-identification," in *ACM MM*, 2018, pp. 274-282.
- [11] Y. Sun, L. Zheng, Y. Yang, Q. Tian, and S. Wang, "Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline)," in *ECCV*, 2018, pp. 480-496.
- [12] L. Zheng, H. Zhang, S. Sun, M. Chandraker, Y. Yang, and Q. Tian, "Person re-identification in the wild," in *CVPR*, 2017, pp. 1367-1376.
- [13] P. Wang, B. Jiao, L. Yang, Y. Yang, S. Zhang, W. Wei, and Y. Zhang, "Vehicle re-identification in aerial imagery: Dataset and approach," in *ICCV*, 2019, pp. 460-469.
- [14] S. Xuan and S. Zhang, "Intra-inter camera similarity for unsupervised person re-identification," in *CVPR*, 2021, pp. 11 926-11 935.
- [15] X. Zhang, D. Li, Z. Wang, J. Wang, E. Ding, J. Q. Shi, Z. Zhang, and J. Wang, "Implicit sample extension for unsupervised person re-identification," in *CVPR*, 2022, pp. 7369-7378.
- [16] H.-X. Yu, W.-S. Zheng, A. Wu, X. Guo, S. Gong, and J.-H. Lai, "Unsupervised person re-identification by soft multilabel learning," in *CVPR*, 2019, pp. 2148-2157.
- [17] M. Ye, W. Ruan, B. Du, and M. Z. Shou, "Channel augmented joint learning for visible-infrared recognition," in *ICCV*, 2021, pp. 13 567-13 576.
- [18] Z. Wu and M. Ye, "Unsupervised visible-infrared person re-identification via progressive graph matching and alternate learning," in *CVPR*, 2023, pp. 9548-9558.
- [19] X. Jin, C. Lan, W. Zeng, Z. Chen, and L. Zhang, "Style normalization and restitution for generalizable person re-identification," in *CVPR*, 2020, pp. 3143-3152.
- [20] Y. Zhao, Z. Zhong, F. Yang, Z. Luo, Y. Lin, S. Li, and N. Sebe, "Learning to generalize unseen domains via memory-based multi-source meta-learning for person re-identification," in *CVPR*, 2021, pp. 6277-6286.
- [21] S. Chen, M. Ye, and B. Du, "Rotation invariant transformer for recognizing object in uavs," in *ACM MM*, 2022, pp. 2565-2574.
- [22] H. Rao, S. Wang, X. Hu, M. Tan, Y. Guo, J. Cheng, X. Liu, and B. Hu, "A self-supervised gait encoding approach with locality-awareness for 3d skeleton based person re-identification," *IEEE TPAMI*, vol. 44, no. 10, pp. 6649-6666, 2021.
- [23] H. Xiao, W. Lin, B. Sheng, K. Lu, J. Yan, J. Wang, E. Ding, Y. Zhang, and H. Xiong, "Group re-identification: Leveraging and integrating multi-grain information," in *ACM MM*, 2018, pp. 192-200.
- [24] L. Wei, S. Zhang, W. Gao, and Q. Tian, "Person transfer gan to bridge domain gap for person re-identification," in *CVPR*, 2018, pp. 79-88.
- [25] S. Li, T. Xiao, H. Li, B. Zhou, D. Yue, and X. Wang, "Person search with natural language description," in *CVPR*, 2017, pp. 1970-1979.
- [26] S. Zhang, Q. Zhang, Y. Yang, X. Wei, P. Wang, B. Jiao, and Y. Zhang, "Person re-identification in aerial imagery," *IEEE TMM*, vol. 23, pp. 281-291, 2020.
- [27] X. Qian, W. Wang, L. Zhang, F. Zhu, Y. Fu, T. Xiang, Y.-G. Jiang, and X. Xue, "Long-term cloth-changing person re-identification," in *ACCV*, 2020.
- [28] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *NeurIPS*, vol. 30, 2017.
- [29] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *ICLR*, 2020.

- [30] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," *arXiv preprint arXiv:2103.14030*, 2021.
- [31] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions," *arXiv preprint arXiv:2102.12122*, 2021.
- [32] S. He, H. Luo, P. Wang, F. Wang, H. Li, and W. Jiang, "Transreid: Transformer-based object re-identification," in *ICCV*, 2021, pp. 15 013–15 022.
- [33] J. Cao, Y. Pang, R. M. Anwer, H. Cholakkal, J. Xie, M. Shah, and F. S. Khan, "Pstr: End-to-end one-step person search with transformers," in *CVPR*, 2022, pp. 9458–9467.
- [34] B. Zhang, Y. Liang, and M. Du, "Interlaced perception for person re-identification based on swin transformer," in *IEEE ICIVC*, 2022, pp. 24–30.
- [35] B. Comandur, "Sports re-id: Improving re-identification of players in broadcast videos of team sports," *arXiv preprint arXiv:2206.02373*, 2022.
- [36] M. Walmer, S. Suri, K. Gupta, and A. Shrivastava, "Teaching matters: Investigating the role of supervision in vision transformers," in *CVPR*, 2023, pp. 7486–7496.
- [37] P. Xu, X. Zhu, and D. A. Clifton, "Multimodal learning with transformers: A survey," *IEEE TPAMI*, 2023.
- [38] K. Han, Y. Wang, H. Chen, X. Chen, J. Guo, Z. Liu, Y. Tang, A. Xiao, C. Xu, Y. Xu *et al.*, "A survey on vision transformer," *IEEE TPAMI*, vol. 45, no. 1, pp. 87–110, 2022.
- [39] M. Naseer, K. Ranasinghe, S. Khan, M. Hayat, F. S. Khan, and M.-H. Yang, "Intriguing properties of vision transformers," *arXiv preprint arXiv:2105.10497*, 2021.
- [40] D. Jiang and M. Ye, "Cross-modal implicit relation reasoning and aligning for text-to-image person retrieval," in *CVPR*, 2023, pp. 2787–2797.
- [41] H. Luo, P. Wang, Y. Xu, F. Ding, Y. Zhou, F. Wang, H. Li, and R. Jin, "Self-supervised pre-training for transformer-based person re-identification," *arXiv preprint arXiv:2111.12084*, 2021.
- [42] H. Rao and C. Miao, "Transg: Transformer-based skeleton graph prototype contrastive learning with structure-trajectory prompted reconstruction for person re-identification," in *CVPR*, 2023, pp. 22 118–22 128.
- [43] H. Li, M. Ye, C. Wang, and B. Du, "Pyramidal transformer with conv-patchify for person re-identification," in *ACM MM*, 2022, pp. 7317–7326.
- [44] S. D. Khan and H. Ullah, "A survey of advances in vision-based vehicle re-identification," *CVIU*, vol. 182, pp. 50–63, 2019.
- [45] Z. Wang, Z. Wang, Y. Zheng, Y. Wu, W. Zeng, and S. Satoh, "Beyond intra-modality: A survey of heterogeneous person re-identification," *arXiv preprint arXiv:1905.10048*, 2019.
- [46] X. Liu, W. Liu, T. Mei, and H. Ma, "A deep learning-based approach to progressive vehicle re-identification for urban surveillance," in *ECCV*. Springer, 2016, pp. 869–884.
- [47] X. Liu, W. Liu, H. Ma, and H. Fu, "Large-scale vehicle re-identification in urban surveillance videos," in *ICME*. IEEE, 2016, pp. 1–6.
- [48] Y.-C. Chen, X. Zhu, W.-S. Zheng, and J.-H. Lai, "Person re-identification by camera correlation aware feature augmentation," *IEEE TPAMI*, vol. 40, no. 2, pp. 392–408, 2017.
- [49] C. Chen, M. Ye, and D. Jiang, "Towards modality-agnostic person re-identification with descriptive query," in *CVPR*, 2023, pp. 15 128–15 137.
- [50] B. Yang, M. Ye, J. Chen, and Z. Wu, "Augmented dual-contrastive aggregation learning for unsupervised visible-infrared person re-identification," in *ACM MM*, 2022, pp. 2843–2851.
- [51] S. Teng, S. Zhang, Q. Huang, and N. Sebe, "Viewpoint and scale consistency reinforcement for uav vehicle re-identification," *IJCV*, vol. 129, pp. 719–735, 2021.
- [52] T. Li, J. Liu, W. Zhang, Y. Ni, W. Wang, and Z. Li, "Uav-human: A large benchmark for human behavior understanding with unmanned aerial vehicles," in *CVPR*, 2021, pp. 16 266–16 275.
- [53] S. Kumar, E. Yaghoubi, A. Das, B. Harish, and H. Proença, "The p-destre: a fully annotated dataset for pedestrian detection, tracking, re-identification and search from aerial devices," *arXiv preprint arXiv:2004.02782*, 2020.
- [54] M. Ye, Z. Wu, C. Chen, and B. Du, "Channel augmentation for visible-infrared re-identification," *IEEE TPAMI*, no. 01, pp. 1–16, 2023.
- [55] Y. Zhai, Y. Zeng, D. Cao, and S. Lu, "Trireid: Towards multi-modal person re-identification via descriptive fusion model," in *ICMR*, 2022, pp. 63–71.
- [56] Z. Wang, Z. Fang, J. Wang, and Y. Yang, "Vita: Visual-textual attributes alignment in person search by natural language," in *ECCV*. Springer, 2020, pp. 402–420.
- [57] M. Li, X. Zhu, and S. Gong, "Unsupervised tracklet person re-identification," *IEEE TPAMI*, vol. 42, no. 7, pp. 1770–1782, 2019.
- [58] D. Fu, D. Chen, J. Bao, H. Yang, L. Yuan, L. Zhang, H. Li, and D. Chen, "Unsupervised pre-training for person re-identification," in *CVPR*, 2021, pp. 14 750–14 759.
- [59] M. Ye, H. Li, B. Du, J. Shen, L. Shao, and S. C. Hoi, "Collaborative refining for person re-identification with label noise," *IEEE TIP*, vol. 31, pp. 379–391, 2021.
- [60] M. Ye, J. Shen, X. Zhang, P. C. Yuen, and S.-F. Chang, "Augmentation invariant and instance spreading feature for softmax embedding," *IEEE TPAMI*, 2020.
- [61] Y. Lin, X. Dong, L. Zheng, Y. Yan, and Y. Yang, "A bottom-up clustering approach to unsupervised person re-identification," in *AAAI*, vol. 33, no. 01, 2019, pp. 8738–8745.
- [62] Y. Cho, W. J. Kim, S. Hong, and S.-E. Yoon, "Part-based pseudo label refinement for unsupervised person re-identification," in *CVPR*, 2022, pp. 7308–7318.
- [63] H. Li, M. Ye, and B. Du, "Weperson: Learning a generalized re-identification model from all-weather virtual data," in *ACM MM*, 2021, pp. 3115–3123.
- [64] Y. Bai, J. Jiao, W. Ce, J. Liu, Y. Lou, X. Feng, and L.-Y. Duan, "Person30k: A dual-meta generalization network for person re-identification," in *CVPR*, 2021, pp. 2123–2132.
- [65] H. Ni, J. Song, X. Luo, F. Zheng, W. Li, and H. T. Shen, "Meta distribution alignment for generalizable person re-identification," in *CVPR*, 2022, pp. 2487–2496.
- [66] L. Fan, T. Li, R. Fang, R. Hristov, Y. Yuan, and D. Katabi, "Learning longterm representations for person re-identification using radio signals," in *CVPR*, 2020, pp. 10 699–10 709.
- [67] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.
- [68] Z. Zheng, L. Zheng, and Y. Yang, "A discriminatively learned cnn embedding for person re-identification," *ACM TOMM*, vol. 14, no. 1, pp. 1–20, 2017.
- [69] H. Luo, W. Jiang, Y. Gu, F. Liu, X. Liao, S. Lai, and J. Gu, "A strong baseline and batch normalization neck for deep person re-identification," *IEEE TMM*, vol. 22, no. 10, pp. 2597–2609, 2019.
- [70] H. Park and B. Ham, "Relation network for person re-identification," in *AAAI*, vol. 34, no. 07, 2020, pp. 11 839–11 847.
- [71] D. Meng, L. Li, X. Liu, Y. Li, S. Yang, Z.-J. Zha, X. Gao, S. Wang, and Q. Huang, "Parsing-based view-aware embedding network for vehicle re-identification," in *CVPR*, 2020, pp. 7103–7112.
- [72] M. M. Kalayeh, E. Basaran, M. Gökmen, M. E. Kamasak, and M. Shah, "Human semantic parsing for person re-identification," in *CVPR*, 2018, pp. 1062–1071.
- [73] Y. Suh, J. Wang, S. Tang, T. Mei, and K. M. Lee, "Part-aligned bilinear representations for person re-identification," in *ECCV*, 2018, pp. 402–419.
- [74] C. Su, J. Li, S. Zhang, J. Xing, W. Gao, and Q. Tian, "Pose-driven deep convolutional model for person re-identification," in *ICCV*, 2017, pp. 3960–3969.
- [75] B. He, J. Li, Y. Zhao, and Y. Tian, "Part-regularized near-duplicate vehicle re-identification," in *CVPR*, 2019, pp. 3997–4005.
- [76] G. Wang, S. Yang, H. Liu, Z. Wang, Y. Yang, S. Wang, G. Yu, E. Zhou, and J. Sun, "High-order information matters: Learning relation and topology for occluded person re-identification," in *CVPR*, 2020, pp. 6449–6458.
- [77] P. Khorramshahi, A. Kumar, N. Peri, S. S. Rambhatla, J.-C. Chen, and R. Chellappa, "A dual-path model with adaptive attention for vehicle re-identification," in *ICCV*, 2019, pp. 6132–6141.
- [78] K. Zhou, Y. Yang, A. Cavallaro, and T. Xiang, "Omni-scale feature learning for person re-identification," in *ICCV*, 2019, pp. 3702–3712.
- [79] H. Guo, K. Zhu, M. Tang, and J. Wang, "Two-level attention network with multi-grain ranking loss for vehicle re-identification," *IEEE TIP*, vol. 28, no. 9, pp. 4328–4338, 2019.
- [80] W. Li, X. Zhu, and S. Gong, "Harmonious attention network for person re-identification," in *CVPR*, 2018, pp. 2285–2294.
- [81] Z. Zhang, C. Lan, W. Zeng, X. Jin, and Z. Chen, "Relation-aware global attention for person re-identification," in *CVPR*, 2020, pp. 3186–3195.

- [82] B. Chen, W. Deng, and J. Hu, "Mixed high-order attention network for person re-identification," in *ICCV*, 2019, pp. 371–381.
- [83] N. McLaughlin, J. M. Del Rincon, and P. Miller, "Recurrent convolutional network for video-based person re-identification," in *CVPR*, 2016, pp. 1325–1334.
- [84] Y. Yan, B. Ni, Z. Song, C. Ma, Y. Yan, and X. Yang, "Person re-identification via recurrent feature aggregation," in *ECCV*. Springer, 2016, pp. 701–716.
- [85] H. Liu, Z. Jie, K. Jayashree, M. Qi, J. Jiang, S. Yan, and J. Feng, "Video-based person re-identification with accumulative motion context," *IEEE transactions on circuits and systems for video technology*, vol. 28, no. 10, pp. 2788–2802, 2017.
- [86] H. Li, C. Li, X. Zhu, A. Zheng, and B. Luo, "Multi-spectral vehicle re-identification: A challenge," in *AAAI*, vol. 34, no. 07, 2020, pp. 11 345–11 353.
- [87] B. Yang, J. Chen, and M. Ye, "Towards grand unified representation learning for unsupervised visible-infrared person re-identification," in *ICCV*, 2023, pp. 11 069–11 079.
- [88] K. Niu, Y. Huang, W. Ouyang, and L. Wang, "Improving description-based person re-identification by multi-granularity image-text alignments," *IEEE TIP*, pp. 5542–5556, 2020.
- [89] Y. Wu, Z. Yan, X. Han, G. Li, C. Zou, and S. Cui, "Lapscore: language-guided person search via color reasoning," in *ICCV*, 2021, pp. 1624–1633.
- [90] Z. Ding, C. Ding, Z. Shao, and D. Tao, "Semantically self-aligned network for text-to-image part-aware person re-identification," *arXiv preprint arXiv:2107.12666*, 2021.
- [91] M. Ye, J. Shen, and L. Shao, "Visible-infrared person re-identification via homogeneous augmented tri-modal learning," *IEEE TIFS*, vol. 16, pp. 728–739, 2020.
- [92] M. Ye, X. Lan, Z. Wang, and P. C. Yuen, "Bi-directional center-constrained top-ranking for visible thermal person re-identification," *IEEE TIFS*, vol. 15, pp. 407–419, 2019.
- [93] M. Ye, X. Lan, J. Li, and P. Yuen, "Hierarchical discriminative learning for visible thermal person re-identification," in *AAAI*, vol. 32, no. 1, 2018.
- [94] S. Zhang, Y. Yang, P. Wang, G. Liang, X. Zhang, and Y. Zhang, "Attend to the difference: Cross-modality person re-identification via contrastive correlation," *IEEE TIP*, vol. 30, pp. 8861–8872, 2021.
- [95] A. Wu, W.-S. Zheng, H.-X. Yu, S. Gong, and J. Lai, "Rgb-infrared cross-modality person re-identification," in *ICCV*, 2017, pp. 5380–5389.
- [96] G. Wang, T. Zhang, J. Cheng, S. Liu, Y. Yang, and Z. Hou, "Rgb-infrared cross-modality person re-identification via joint pixel and feature alignment," in *ICCV*, 2019, pp. 3623–3632.
- [97] Z. Wang, Z. Wang, Y. Zheng, Y.-Y. Chuang, and S. Satoh, "Learning to reduce dual-level discrepancy for infrared-visible person re-identification," in *CVPR*, 2019, pp. 618–626.
- [98] G.-A. Wang, T. Zhang, Y. Yang, J. Cheng, J. Chang, X. Liang, and Z.-G. Hou, "Cross-modality paired-images generation for rgb-infrared person re-identification," in *AAAI*, vol. 34, no. 07, 2020, pp. 12 144–12 151.
- [99] Y. Zhang and H. Lu, "Deep cross-modal projection learning for image-text matching," in *ECCV*, 2018, pp. 686–701.
- [100] N. Sarafianos, X. Xu, and I. A. Kakadiaris, "Adversarial representation learning for text-to-image matching," in *ICCV*, 2019, pp. 5814–5824.
- [101] Z. Shao, X. Zhang, M. Fang, Z. Lin, J. Wang, and C. Ding, "Learning granularity-unified representations for text-to-image person re-identification," in *ACM MM*, 2022, pp. 5566–5574.
- [102] A. Farooq, M. Awais, J. Kittler, and S. S. Khalid, "Axm-net: Implicit cross-modal feature alignment for person re-identification," in *AAAI*, vol. 36, no. 4, 2022, pp. 4477–4485.
- [103] Y. Ge, F. Zhu, D. Chen, R. Zhao *et al.*, "Self-paced contrastive learning with hybrid memory for domain adaptive object re-id," *NeurIPS*, vol. 33, pp. 11 309–11 321, 2020.
- [104] Y. Dai, J. Liu, Y. Sun, Z. Tong, C. Zhang, and L.-Y. Duan, "Idm: An intermediate domain module for domain adaptive person re-id," in *ICCV*, 2021, pp. 11 864–11 874.
- [105] Z. Bai, Z. Wang, J. Wang, D. Hu, and E. Ding, "Unsupervised multi-source domain adaptation for person re-identification," in *CVPR*, 2021, pp. 12 914–12 923.
- [106] K. Zheng, W. Liu, L. He, T. Mei, J. Luo, and Z.-J. Zha, "Group-aware label transfer for domain adaptive person re-identification," in *CVPR*, 2021, pp. 5310–5319.
- [107] Y. Lin, L. Xie, Y. Wu, C. Yan, and Q. Tian, "Unsupervised person re-identification via softened similarity learning," in *CVPR*, 2020, pp. 3390–3399.
- [108] D. Wang and S. Zhang, "Unsupervised person re-identification via multi-label classification," in *CVPR*, 2020, pp. 10 981–10 990.
- [109] W. Deng, L. Zheng, Q. Ye, G. Kang, Y. Yang, and J. Jiao, "Image-image domain adaptation with preserved self-similarity and domain-dissimilarity for person re-identification," in *CVPR*, 2018, pp. 994–1003.
- [110] X. Zhang, Y. Ge, Y. Qiao, and H. Li, "Refining pseudo labels with clustering consensus over generations for unsupervised object re-identification," in *CVPR*, 2021, pp. 3436–3445.
- [111] L. Wu, D. Liu, W. Zhang, D. Chen, Z. Ge, F. Boussaid, M. Benamoun, and J. Shen, "Pseudo-pair based self-similarity learning for unsupervised person re-identification," *IEEE TIP*, vol. 31, pp. 4803–4816, 2022.
- [112] Z. Dai, G. Wang, W. Yuan, S. Zhu, and P. Tan, "Cluster contrast for unsupervised person re-identification," in *ACCV*, 2022, pp. 1142–1160.
- [113] Q. Yang, A. Wu, and W.-S. Zheng, "Person re-identification by contour sketch under moderate clothing change," *IEEE TPAMI*, vol. 43, no. 6, pp. 2029–2046, 2019.
- [114] X. Jia, X. Zhong, M. Ye, W. Liu, and W. Huang, "Complementary data augmentation for cloth-changing person re-identification," *IEEE TIP*, vol. 31, pp. 4227–4239, 2022.
- [115] W. Xu, H. Liu, W. Shi, Z. Miao, Z. Lu, and F. Chen, "Adversarial feature disentanglement for long-term person re-identification," in *IJCAI*, 2021, pp. 1201–1207.
- [116] F. Liu, M. Ye, and B. Du, "Dual level adaptive weighting for cloth-changing person re-identification," *IEEE TIP*, 2023.
- [117] P. Guo, H. Liu, J. Wu, G. Wang, and T. Wang, "Semantic-aware consistency network for cloth-changing person re-identification," *arXiv preprint arXiv:2308.14113*, 2023.
- [118] X. Jin, T. He, K. Zheng, Z. Yin, X. Shen, Z. Huang, R. Feng, J. Huang, Z. Chen, and X.-S. Hua, "Cloth-changing person re-identification from a single image with gait prediction and regularization," in *CVPR*, 2022, pp. 14 278–14 287.
- [119] P. Hong, T. Wu, A. Wu, X. Han, and W.-S. Zheng, "Fine-grained shape-appearance mutual learning for cloth-changing person re-identification," in *CVPR*, 2021, pp. 10 513–10 522.
- [120] S. Choi, T. Kim, M. Jeong, H. Park, and C. Kim, "Meta batch-instance normalization for generalizable person re-identification," in *CVPR*, 2021, pp. 3425–3435.
- [121] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [122] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, "Masked autoencoders are scalable vision learners," in *CVPR*, 2022, pp. 16 000–16 009.
- [123] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, "Emerging properties in self-supervised vision transformers," in *ICCV*, 2021, pp. 9650–9660.
- [124] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *NeurIPS*, vol. 33, pp. 1877–1901, 2020.
- [125] M. Dehghani, J. Djolonga, B. Mustafa, P. Padlewski, J. Heek, J. Gilmer, A. P. Steiner, M. Caron, R. Geirhos, I. Alabdulmohsin *et al.*, "Scaling vision transformers to 22 billion parameters," in *ICML*. PMLR, 2023, pp. 7480–7512.
- [126] X. Zhai, A. Kolesnikov, N. Houlsby, and L. Beyer, "Scaling vision transformers," in *CVPR*, 2022, pp. 12 104–12 113.
- [127] G. Zhang, P. Zhang, J. Qi, and H. Lu, "Hat: Hierarchical aggregation transformers for person re-identification," in *ACM MM*, 2021, pp. 516–525.
- [128] F. Shen, Y. Xie, J. Zhu, X. Zhu, and H. Zeng, "Git: Graph interactive transformer for vehicle re-identification," *IEEE TIP*, vol. 32, pp. 1039–1051, 2023.
- [129] W. Li, C. Zou, M. Wang, F. Xu, J. Zhao, R. Zheng, Y. Cheng, and W. Chu, "Dc-former: Diverse and compact transformer for person re-identification," *arXiv preprint arXiv:2302.14335*, 2023.
- [130] B. Tan, L. Xu, Z. Qiu, Q. Wu, and F. Meng, "Mfat: A multi-level feature aggregated transformer for person re-identification," in *ICASSP*. IEEE, 2023, pp. 1–5.
- [131] X. Chen, C. Xu, Q. Cao, J. Xu, Y. Zhong, J. Xu, Z. Li, J. Wang, and S. Gao, "Oh-former: Omni-relational high-order transformer for person re-identification," *arXiv preprint arXiv:2109.11159*, 2021.

- [132] H. Zhu, W. Ke, D. Li, J. Liu, L. Tian, and Y. Shan, "Dual cross-attention learning for fine-grained visual categorization and object re-identification," in *CVPR*, 2022, pp. 4692–4702.
- [133] X. Tian, J. Liu, Z. Zhang, C. Wang, Y. Qu, Y. Xie, and L. Ma, "Hierarchical walking transformer for object re-identification," in *ACM MM*, 2022, pp. 4224–4232.
- [134] L. Shen, T. He, Y. Guo, and G. Ding, "X-reid: Cross-instance transformer for identity-level person re-identification," *arXiv preprint arXiv:2302.02075*, 2023.
- [135] J. Mao, Y. Yao, Z. Sun, X. Huang, F. Shen, and H.-T. Shen, "Attention map guided transformer pruning for occluded person re-identification on edge device," *IEEE TMM*, 2023.
- [136] K. Zhu, H. Guo, S. Zhang, Y. Wang, G. Huang, H. Qiao, J. Liu, J. Wang, and M. Tang, "Aaformer: Auto-aligned transformer for person re-identification," *arXiv preprint arXiv:2104.00921*, 2021.
- [137] S. Lai, Z. Chai, and X. Wei, "Transformer meets part model: Adaptive part division for person re-identification," in *ICCV*, 2021, pp. 4150–4157.
- [138] Y. Li, J. He, T. Zhang, X. Liu, Y. Zhang, and F. Wu, "Diverse part discovery: Occluded person re-identification with part-aware transformer," in *CVPR*, 2021, pp. 2898–2907.
- [139] T. Wang, H. Liu, P. Song, T. Guo, and W. Shi, "Pose-guided feature disentangling for occluded person re-identification based on transformer," in *AAAI*, vol. 36, no. 3, 2022, pp. 2540–2549.
- [140] M. Zhou, H. Liu, Z. Lv, W. Hong, and X. Chen, "Motion-aware transformer for occluded person re-identification," *arXiv preprint arXiv:2202.04243*, 2022.
- [141] C. Chen, M. Ye, M. Qi, and B. Du, "Sketch transformer: Asymmetrical disentanglement learning from dynamic synthesis," in *ACM MM*, 2022, pp. 4012–4020.
- [142] W. Qian, H. Luo, S. Peng, F. Wang, C. Chen, and H. Li, "Unstructured feature decoupling for vehicle re-identification," in *ECCV*. Springer, 2022, pp. 336–353.
- [143] S. G. Mallat, "A theory for multiresolution signal decomposition: the wavelet representation," *IEEE TPAMI*, vol. 11, no. 7, pp. 674–693, 1989.
- [144] G. Zhang, Y. Zhang, T. Zhang, B. Li, and S. Pu, "Pha: Patch-wise high-frequency augmentation for transformer-based person re-identification," in *CVPR*, 2023, pp. 14 133–14 142.
- [145] M. Jia, X. Cheng, S. Lu, and J. Zhang, "Learning disentangled representation implicitly via transformer for occluded person re-identification," *IEEE TMM*, vol. 25, pp. 1294–1305, 2022.
- [146] B. Xu, L. He, J. Liang, and Z. Sun, "Learning feature recovery transformer for occluded person re-identification," *IEEE TIP*, vol. 31, pp. 4651–4662, 2022.
- [147] Y. Ye, H. Zhou, J. Yu, Q. Hu, and W. Yang, "Dynamic feature pruning and consolidation for occluded person re-identification," *arXiv preprint arXiv:2211.14742*, 2022.
- [148] T. Wang, H. Liu, W. Li, M. Ban, T. Guo, and Y. Li, "Feature completion transformer for occluded person re-identification," *arXiv preprint arXiv:2303.01656*, 2023.
- [149] X. Cheng, M. Jia, Q. Wang, and J. Zhang, "More is better: Multi-source dynamic parsing attention for occluded person re-identification," in *ACM MM*, 2022, pp. 6840–6849.
- [150] H. Wang, J. Shen, Y. Liu, Y. Gao, and E. Gavves, "Nformer: Robust person re-identification with neighbor transformer," in *CVPR*, 2022, pp. 7297–7307.
- [151] T. He, X. Jin, X. Shen, J. Huang, Z. Chen, and X.-S. Hua, "Dense interaction learning for video-based person re-identification," in *ICCV*, 2021, pp. 1490–1501.
- [152] X. Liu, C. Yu, P. Zhang, and H. Lu, "Deeply coupled convolution-transformer with spatial-temporal complementary learning for video-based person re-identification," *IEEE TNNLS*, 2023.
- [153] Z. Tang, R. Zhang, Z. Peng, J. Chen, and L. Lin, "Multi-stage spatio-temporal aggregation transformer for video person re-identification," *IEEE TMM*, 2022.
- [154] J. Wu, L. He, W. Liu, Y. Yang, Z. Lei, T. Mei, and S. Z. Li, "Cavit: Contextual alignment vision transformer for video object re-identification," in *ECCV*. Springer, 2022, pp. 549–566.
- [155] X. Liu, P. Zhang, C. Yu, H. Lu, X. Qian, and X. Yang, "A video is worth three views: Trigeminal transformers for video-based person re-identification," *arXiv preprint arXiv:2104.01745*, 2021.
- [156] T. Zhang, L. Wei, L. Xie, Z. Zhuang, Y. Zhang, B. Li, and Q. Tian, "Spatiotemporal transformer for video-based person re-identification," *arXiv preprint arXiv:2103.16469*, 2021.
- [157] H. Li, G. Wu, and W.-S. Zheng, "Combined depth space based architecture search for person re-identification," in *CVPR*, 2021, pp. 6729–6738.
- [158] H. Chen, B. Lagadec, and F. Bremond, "Ice: Inter-instance contrastive encoding for unsupervised person re-identification," in *ICCV*, 2021, pp. 14 960–14 969.
- [159] K. Zhu, H. Guo, T. Yan, Y. Zhu, J. Wang, and M. Tang, "Pass: Part-aware self-supervised pre-training for person re-identification," in *ECCV*. Springer Nature Switzerland Cham, 2022, pp. 198–214.
- [160] Y. Wang, G. Qi, S. Li, Y. Chai, and H. Li, "Body part-level domain alignment for domain-adaptive person re-identification with transformer framework," *IEEE TIFS*, vol. 17, pp. 3321–3334, 2022.
- [161] R. Wei, J. Gu, S. He, and W. Jiang, "Transformer-based domain-specific representation for unsupervised domain adaptive vehicle re-identification," *IEEE TITS*, vol. 24, no. 3, pp. 2935–2946, 2022.
- [162] S. Liao and L. Shao, "Transmatcher: Deep image matching through transformers for generalizable person re-identification," *NeurIPS*, vol. 34, pp. 1992–2003, 2021.
- [163] H. Ni, Y. Li, L. Gao, H. T. Shen, and J. Song, "Part-aware transformer for generalizable person re-identification," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 11 280–11 289.
- [164] X. Wang, X. Wang, B. Jiang, and B. Luo, "Few-shot learning meets transformer: Unified query-support transformers for few-shot classification," *IEEE TCSVT*, 2023.
- [165] M. Ye, Y. Cheng, X. Lan, and H. Zhu, "Improving night-time pedestrian retrieval with distribution alignment and contextual distance," *IEEE TII*, vol. 16, no. 1, pp. 615–624, 2019.
- [166] K. Jiang, T. Zhang, X. Liu, B. Qian, Y. Zhang, and F. Wu, "Cross-modality transformer for visible-infrared person re-identification," in *ECCV*. Springer, 2022, pp. 480–496.
- [167] T. Liang, Y. Jin, W. Liu, and Y. Li, "Cross-modality transformer with modality mining for visible-infrared person re-identification," *IEEE TMM*, 2023.
- [168] C. Chen, M. Ye, M. Qi, J. Wu, J. Jiang, and C.-W. Lin, "Structure-aware positional transformer for visible-infrared person re-identification," *IEEE TIP*, vol. 31, pp. 2352–2364, 2022.
- [169] Y. Feng, J. Yu, F. Chen, Y. Ji, F. Wu, S. Liu, and X.-Y. Jing, "Visible-infrared person re-identification via cross-modality interaction transformer," *IEEE TMM*, 2022.
- [170] B. Yang, J. Chen, and M. Ye, "Top-k visual tokens transformer: Selecting tokens for visible-infrared person re-identification," in *ICASSP*. IEEE, 2023, pp. 1–5.
- [171] J. Zhao, H. Wang, Y. Zhou, R. Yao, S. Chen, and A. El Saddik, "Spatial-channel enhanced transformer for visible-infrared person re-identification," *IEEE TMM*, 2022.
- [172] M. Ye, C. Liang, Z. Wang, Q. Leng, J. Chen, and J. Liu, "Specific person retrieval via incomplete text description," in *ACM ICMRI*, 2015, pp. 547–550.
- [173] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *ICML*. PMLR, 2021, pp. 8748–8763.
- [174] X. Han, S. He, L. Zhang, and T. Xiang, "Text-based person search with limited data," 2021.
- [175] S. Yan, N. Dong, L. Zhang, and J. Tang, "Clip-driven fine-grained text-image person re-identification," *arXiv preprint arXiv:2210.10276*, 2022.
- [176] J. Zuo, C. Yu, N. Sang, and C. Gao, "Plip: Language-image pre-training for person representation learning," *arXiv preprint arXiv:2305.08386*, 2023.
- [177] S. Li, L. Sun, and Q. Li, "Clip-reid: Exploiting vision-language model for image re-identification without concrete text labels," in *AAAI*, vol. 37, no. 1, 2023, pp. 1405–1413.
- [178] G. Wang, F. Yu, J. Li, Q. Jia, and S. Ding, "Exploiting the textual potential from vision-language pre-training for text-based person search," *arXiv preprint arXiv:2303.04497*, 2023.
- [179] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, "Learning to prompt for vision-language models," *IJCV*, vol. 130, no. 9, pp. 2337–2348, 2022.
- [180] Y. Zhang, Y. Wang, H. Li, and S. Li, "Cross-compatible embedding and semantic consistent feature construction for sketch re-identification," in *ACM MM*, 2022, pp. 3347–3355.

- [181] K. W. Lee, B. Jawade, D. Mohan, S. Setlur, and V. Govindaraju, "Attribute de-biased vision transformer (ad-vit) for long-term person re-identification," in *IEEE AVSS*. IEEE, 2022, pp. 1–8.
- [182] Y. Ci, Y. Wang, M. Chen, S. Tang, L. Bai, F. Zhu, R. Zhao, F. Yu, D. Qi, and W. Ouyang, "Unihcp: A unified model for human-centric perceptions," in *CVPR*, 2023, pp. 17 840–17 852.
- [183] S. Tang, C. Chen, Q. Xie, M. Chen, Y. Wang, Y. Ci, L. Bai, F. Zhu, H. Yang, L. Yi *et al.*, "Humanbench: Towards general human-centric perception with projector assisted pretraining," in *CVPR*, 2023, pp. 21 970–21 982.
- [184] W. Chen, X. Xu, J. Jia, H. Luo, Y. Wang, F. Wang, R. Jin, and X. Sun, "Beyond appearance: a semantic controllable self-supervised learning framework for human-centric visual tasks," in *CVPR*, 2023, pp. 15 050–15 061.
- [185] R. Yu, D. Du, R. LaLonde, D. Davila, C. Funk, A. Hoogs, and B. Clipp, "Cascade transformers for end-to-end person search," in *CVPR*, 2022, pp. 7267–7276.
- [186] W. Zheng, S. Gong, and T. Xiang, "Associating groups of people," in *BMVC*, 2009, pp. 1–11.
- [187] Q. Zhang, J.-H. Lai, Z. Feng, and X. Xie, "Uncertainty modeling with second-order transformer for group re-identification," in *AAAI*, vol. 36, no. 3, 2022, pp. 3318–3325.
- [188] D. Organisciak, M. Poyser, A. Alsehim, S. Hu, B. K. Isaac-Medina, T. P. Breckon, and H. P. Shum, "Uav-reid: A benchmark on unmanned aerial vehicle re-identification in video imagery," *arXiv preprint arXiv:2104.06219*, 2021.
- [189] S. N. Ferdous, X. Li, and S. Lyu, "Uncertainty aware multitask pyramid vision transformer for uav-based object re-identification," in *ICIP*. IEEE, 2022, pp. 2381–2385.
- [190] Z. Xie, Z. Zhang, Y. Cao, Y. Lin, J. Bao, Z. Yao, Q. Dai, and H. Hu, "Simmim: A simple framework for masked image modeling," in *CVPR*, 2022, pp. 9653–9663.
- [191] D. T. Nguyen, H. G. Hong, K. W. Kim, and K. R. Park, "Person recognition system based on a combination of body images from visible light and thermal cameras," *Sensors*, vol. 17, no. 3, p. 605, 2017.
- [192] J. Wang, Z. Zhang, M. Chen, Y. Zhang, C. Wang, B. Sheng, Y. Qu, and Y. Xie, "Optimal transport for label-efficient visible-infrared person re-identification," in *ECCV*. Springer, 2022, pp. 93–109.
- [193] M. Korschens and J. Denzler, "Elpephants: A fine-grained dataset for elephant re-identification," in *ICCV Workshop*, 2019, pp. 0–0.
- [194] L. Wang, R. Ding, Y. Zhai, Q. Zhang, W. Tang, N. Zheng, and G. Hua, "Giant panda identification," *IEEE TIP*, vol. 30, pp. 2837–2849, 2021.
- [195] A. Howard, i. Ken, Southerland, R. Holbrook, and T. Cheeseman, "Happywhale - whale and dolphin identification," 2022. [Online]. Available: <https://kaggle.com/competitions/happy-whale-and-dolphin>
- [196] L. I. Kuncheva, F. Williams, S. L. Hennessey, and J. J. Rodríguez, "A benchmark database for animal re-identification and tracking," in *IEEE IPAS*. IEEE, 2022, pp. 1–6.
- [197] S. Schneider, G. W. Taylor, S. Linquist, and S. C. Kremer, "Past, present and future approaches using computer vision for animal re-identification from camera trap data," *Methods in Ecology and Evolution*, vol. 10, no. 4, pp. 461–470, 2019.
- [198] T. Zhang, Q. Zhao, C. Da, L. Zhou, L. Li, and S. Jiancui, "Yakreid-103: A benchmark for yak re-identification," in *IEEE IJCB*. IEEE, 2021, pp. 1–8.
- [199] E. Nepovninnykh, T. Eerola, V. Biard, P. Mutka, M. Niemi, M. Kunasranta, and H. Kälviäinen, "Sealid: Saimaa ringed seal re-identification dataset," *Sensors*, vol. 22, no. 19, p. 7602, 2022.
- [200] K. Papafitsoros, L. Adam, V. Čermák, and L. Pícek, "Seaturtleid: A novel long-span dataset highlighting the importance of timestamps in wildlife re-identification," *arXiv preprint arXiv:2211.10307*, 2022.
- [201] L. Bergamini, A. Porrello, A. C. Dondona, E. Del Negro, M. Mattioli, N. D'alterio, and S. Calderara, "Multi-views embedding for cattle re-identification," in *IEEE SITIS*, 2018, pp. 184–191.
- [202] S. Bouma, M. D. Pawley, K. Hupman, and A. Gilman, "Individual common dolphin identification via metric embedding learning," in *IEEE IVCNZ*, 2018, pp. 1–6.
- [203] J. Bruslund Haurum, A. Karpova, M. Pedersen, S. Hein Bengtson, and T. B. Moeslund, "Re-identification of zebrafish using metric learning," in *WACV Workshop*, 2020, pp. 1–11.
- [204] S. Li, L. Fu, Y. Sun, Y. Mu, L. Chen, J. Li, and H. Gong, "Individual dairy cow identification based on lightweight convolutional neural network," *Plos one*, vol. 16, no. 11, p. e0260510, 2021.
- [205] J. Gao, T. Burghardt, W. Andrew, A. W. Dowsey, and N. W. Campbell, "Towards self-supervision for video identification of individual holstein-friesian cattle: The cows2021 dataset," *arXiv preprint arXiv:2105.01938*, 2021.
- [206] "Beluga id 2022." [Online]. Available: <https://lila.science/datasets/beluga-id-2022/>
- [207] J. Chan, H. Carrión, R. Mégret, J. L. Agosto-Rivera, and T. Giray, "Honeybee re-identification in video: New datasets and impact of self-supervision." in *VISIGRAPP (5: VISAPP)*, 2022, pp. 517–525.
- [208] "Leopard id 2022." [Online]. Available: <https://lila.science/datasets/leopard-id-2022/>
- [209] "Hyena id 2022." [Online]. Available: <https://lila.science/datasets/hyena-id-2022/>
- [210] M. Zuerl, R. Dirauf, F. Koeferl, N. Steinlein, J. Sueskind, D. Zanca, I. Brehm, L. v. Fersen, and B. Eskofier, "Polarbearvidid: A video-based re-identification benchmark dataset for polar bears," *Animals*, vol. 13, no. 5, p. 801, 2023.
- [211] B. Jiao, L. Liu, L. Gao, R. Wu, G. Lin, P. Wang, and Y. Zhang, "Toward re-identifying any animal," in *NeurIPS*, 2023.
- [212] S. Liao, Y. Hu, X. Zhu, and S. Z. Li, "Person re-identification by local maximal occurrence representation and metric learning," in *CVPR*, 2015, pp. 2197–2206.
- [213] O. Moskvyyak, F. Maire, F. Dayoub, A. O. Armstrong, and M. Baktashmotlagh, "Robust re-identification of manta rays from natural markings by learning pose invariant embeddings," in *DICTA*. IEEE, 2021, pp. 1–8.
- [214] A. Porrello, L. Bergamini, and S. Calderara, "Robust re-identification by multiple views knowledge distillation," in *ECCV*. Springer, 2020, pp. 93–110.
- [215] H. Weideman, C. Stewart, J. Parham, J. Holmberg, K. Flynn, J. Calambokidis, D. B. Paul, A. Bedetti, M. Henley, F. Pope *et al.*, "Extracting identifying contours for african elephants and humpback whales using a learned appearance model," in *WACV*, 2020, pp. 1276–1285.
- [216] T. Cheeseman, K. Southerland, J. Park, M. Olio, K. Flynn, J. Calambokidis, L. Jones, C. Garrigue, A. Frisch Jordan, A. Howard *et al.*, "Advanced image recognition: a fully automated, high-accuracy photo-identification matching system for humpback whales," *Mammalian Biology*, vol. 102, no. 3, pp. 915–929, 2022.
- [217] H. J. Weideman, Z. M. Jablons, J. Holmberg, K. Flynn, J. Calambokidis, R. B. Tyson, J. B. Allen, R. S. Wells, K. Hupman, K. Urian *et al.*, "Integral curvature representation and matching algorithms for identification of dolphins and whales," in *ICCV Workshop*, 2017, pp. 2831–2839.
- [218] D. A. Konovalov, S. Hillcoat, G. Williams, R. A. Birtles, N. Gardiner, and M. I. Curnock, "Individual minke whale recognition using deep learning convolutional neural networks," *Journal of Geoscience and Environment Protection*, vol. 6, pp. 25–36, 2018.
- [219] O. Moskvyyak, F. Maire, F. Dayoub, and M. Baktashmotlagh, "Learning landmark guided embeddings for animal re-identification," in *WACV Workshop*, 2020, pp. 12–19.
- [220] E. Nepovninnykh, T. Eerola, and H. Kälviäinen, "Siamese network based pelage pattern matching for ringed seal re-identification," in *WACV Workshop*, 2020, pp. 25–34.
- [221] G. Koch, R. Zemel, R. Salakhutdinov *et al.*, "Siamese neural networks for one-shot image recognition," in *ICML workshop*, vol. 2, no. 1. Lille, 2015.
- [222] J. Parham, J. Crall, C. Stewart, T. Berger-Wolf, and D. I. Rubenstein, "Animal population censusing at scale with citizen science and photographic identification," in *AAAI*, 2017.
- [223] Y. Zhang, K. Gong, K. Zhang, H. Li, Y. Qiao, W. Ouyang, and X. Yue, "Meta-transformer: A unified framework for multimodal learning," *arXiv preprint arXiv:2307.10802*, 2023.
- [224] N. Ding, Y. Qin, G. Yang, F. Wei, Z. Yang, Y. Su, S. Hu, Y. Chen, C.-M. Chan, W. Chen *et al.*, "Parameter-efficient fine-tuning of large-scale pre-trained language models," *Nature Machine Intelligence*, vol. 5, no. 3, pp. 220–235, 2023.
- [225] N. Pu, Z. Zhong, N. Sebe, and M. S. Lew, "A memorizing and generalizing framework for lifelong person re-identification," *IEEE TPAMI*, 2023.
- [226] J. Gu, H. Luo, K. Wang, W. Jiang, Y. You, and J. Zhao, "Color prompting for data-free continual unsupervised domain adaptive person re-identification," *arXiv preprint arXiv:2308.10716*, 2023.