

Not Every Patch is Needed: Towards a More Efficient and Effective Backbone for Video-based Person Re-identification

Lanyun Zhu, Tianrun Chen, Deyi Ji, Jieping Ye, *Fellow, IEEE*, and Jun Liu

Abstract—This paper proposes a new effective and efficient plug-and-play backbone for video-based person re-identification (ReID). Conventional video-based ReID methods typically use CNN or transformer backbones to extract deep features for every position in every sampled video frame. Here, we argue that this exhaustive feature extraction could be unnecessary, since we find that different frames in a ReID video often exhibit small differences and contain many similar regions due to the relatively slight movements of human beings. Inspired by this, a more selective, efficient paradigm is explored in this paper. Specifically, we introduce a patch selection mechanism to reduce computational cost by choosing only the crucial and non-repetitive patches for feature extraction. Additionally, we present a novel network structure that generates and utilizes pseudo frame global context to address the issue of incomplete views resulting from sparse inputs. By incorporating these new designs, our backbone can achieve both high performance and low computational cost. Extensive experiments on multiple datasets show that our approach reduces the computational cost by 74% compared to ViT-B and 28% compared to ResNet50, while the accuracy is on par with ViT-B and outperforms ResNet50 significantly.

Index Terms—Video-based person re-identification, efficient network.

I. INTRODUCTION

PERSON ReID [91], [37], [9], [92], [58], [93], [30], [31], [71] is an important research topic that has been actively investigated for decades. In recent years, for the industry and research community, video-based ReID [35], [87], [26], [49], [1], [18], [6], [44] is catching more attention with its advantageous properties like containing additional temporal cues and more comprehensive appearance information compared to image-based methods. To unleash the potential of video-based ReID, previous works have explored various techniques to extract effective features from videos. For example, [36], [22], [86] exploit temporal information to enhance person

representations; [76], [73], [75] utilize spatial-temporal correlations based on graph neural networks. The recent success of vision transformers in computer vision has propelled the development in this domain, opening up possibilities for further improvements towards new performance heights.

Despite their success, these existing video-based ReID methods still suffer from a common problem – huge computational burden. Most of the above-mentioned approaches process the video by passing each of the sampled frames through a backbone network to get a frame feature and then merge different frames to form a representation of the entire video. Although this is a typical way of processing video data for deep networks, it will cause the computation cost to rise linearly with the number of sampled frames in a video clip, thus resulting in a huge computational load. This challenge is even more severe when using the latest transformer-based backbone networks like ViT, in which videos are processed into multiple patches and a huge computational cost is required to correlate them. Consequently, although with better performance than CNN-based methods, it becomes almost impossible to use these transformer backbones in practical applications for ReID, especially in many real-world cases where real-time processing is required.

This work aims to address this challenge for computation by providing an innovative alternative path. We begin by asking a more fundamental question: *Is it really necessary to put so many computational resources just to process every region in a frame to fully exploit the information contained in a video?* After analyzing videos used for person ReID in detail, we derive a ‘No’ answer to this question. Specifically, we found two properties in the videos used for ReID that can motivate our computational-efficient solutions. Firstly, the region of interest (the human) can remain in a relatively consistent position in all frames, since the videos are often pre-processed and cropped around target persons. Secondly, the variation caused by human motion is also minimal because intense movements like jumping or striking rarely appear in ReID videos. Considering the two properties, we find that ReID videos typically exhibit minimal temporal variations and have a significant number of similar regions existing in different frames. Intuitively, the repetitive feature extraction for these similar regions is redundant and unnecessary.

Based on the above insights, in this paper, we present a novel and plug-and-play backbone framework for video ReID that leverages the performance advantage of ViT but in a more efficient way with less computational burden. Our proposed

Lanyun Zhu is with Information Systems Technology and Design Pillar, Singapore University of Technology and Design, Singapore, 487372 (e-mail: lanyun_zhu@mymail.sutd.edu.sg)

Tianrun Chen is with College of Computer Science and Technology, Zhejiang University, China, 310027 (email: tianrun.chen@zju.edu.cn)

Deyi Ji and Jieping Ye are with Alibaba Group, China, 310023 (email: jideyi.jdy@alibaba-inc.com, yejieping.ye@alibaba-inc.com)

Jun Liu is with School of Computing and Communications, Lancaster University, UK (e-mail: j.liu81@lancaster.ac.uk)

This work is supported by the National Research Foundation, Singapore under its AI Singapore Programme (AISG Award No: AISG2-PhD-2021-08-006, and AISG-100E-2023-121).

Corresponding author: Jun Liu.

framework achieves excellent performance while maintaining reasonable computation costs, achieving competitive results as ViT but only requiring even lower computations than existing CNN-based backbones. The balance of accuracy and efficiency is achieved by our central idea of pruning redundant patches within each frame and selecting only the crucial patches as transformer input for deep feature extraction. By doing so, we effectively reduce both computation requirement and feature redundancy, thus resulting in a combination of ultra-performance and efficiency. Specifically, we first propose a novel mechanism for patch selection. Differing from previous methods of token pruning [7], [65] in other tasks, our patch selection is carefully designed with a frame-progressive strategy and multi-level selection-determinate features that are tailored for the ReID task. By using it, crucial patches that are useful for ReID can be automatically selected while redundant ones that repeat across frames or exhibit unnecessary properties can be pruned. Next, we design a new feature extractor, namely patch-sparse transformer (PSFormer) to better utilize the selected patches. The motivation of this new extractor is that we find conventional feature extractors suffer from a loss of global information when handling the sparse input from each frame and thus yielding suboptimal results. Therefore, we propose a self-adaptive dynamic routing mechanism to generate pseudo frame global context in our PSFormer, making it capable of effectively handling sparse inputs.

In terms of extracting cross-frame temporal correlations, which is important in video-based ReID, we also propose a novel and efficient approach by using Group of Pictures (GOP). We noted that earlier works explored methods like RNNs [86], [55], optical flows [5], [45], GNNs [76], [75], [73], or events [3] to extract these correlations, but they typically require high computation costs or additional sensors, which is contradictory to our high-efficiency goal. We hereby take an alternative path by pioneering introducing the use of GOP to represent temporal correlations, which is a data structure encoded in compressed videos. For real-world ReID systems, video compression is a necessary step for data transmission between camera terminals and servers. Therefore, the GOP can be directly obtained from the data received by the server, without additional computational costs. Furthermore, the GOP stores information patch-by-patch, which is inherently aligned with the data format of the transformers. Therefore, we use GOP in both the patch selection mechanism and PSFormer in our approach, resulting in a very efficient ReID framework.

We conduct extensive experiments on multiple datasets and the results show the high effectiveness and efficiency of our method. Specifically, our backbone reduces the computational cost by 74% compared to ViT-B and 28% compared to ResNet50, while the results are on par with ViT-B and outperform ResNet50 significantly. Our method is also tested as a plug-and-play backbone for existing methods, with excellent results demonstrating its high generality. In summary, the main contributions of this work are as follows:

- We propose a patch selection mechanism to avoid redundant feature extraction on repetitive or unimportant regions, thereby reducing computational overhead.
- We introduce a patch-sparse transformer to address the

issue of context loss when dealing with sparse inputs, thereby enhancing the model’s effectiveness.

- By incorporating the patch selection mechanism and patch-sparse transformer, we get a novel backbone that can simultaneously achieve high effectiveness and efficiency in video ReID, as demonstrated by our extensive experiments on multiple datasets.

II. RELATED WORK

A. Video-based Person ReID

Video-based person ReID aims to extract comprehensive person representations from videos to facilitate accurate identification. Previous methods have explored mechanisms like RNNs [86], [55], 3D CNNs [36], [22], [46], optical flows [5], [45], events [3], attentions [77], [66], and GNNs [76], [75], [73] to capture cross-frame temporal information. For example, [77] proposes a feature space projection module to preserve more frequency information and avoid feature fragmentation, and introduces a global low-frequency enhancement module to capture global low-frequency features and establish spatial-temporal relationships across video sequences; [36] proposes a two-stream convolution network built from 3D CNNs to capture both spatial and temporal information to help enhance video ReID; [5] improves video ReID by introducing a competitive snippet-similarity aggregation method with the help of optical flow information; [3] utilizes information extracted from event cameras to address the inevitable degradation issues such as motion blur that can adversely affect video ReID. [47] regards the complex temporal features of video ReID as a kind of signal and converts it into frequency domain to extract the useful spectrum cues. [80] introduces a temporal memory diffusion module after the Clip model to extract temporal information by leveraging the sequence-level relations. However, these methods typically require building additional modules that incur extra computational costs to extract temporal cues. In contrast, we propose the first method that uses compressed videos to extract temporal features, which is more efficient and does not require extra sensors or computations like the above-mentioned previous methods. Other methods [20], [61], [89], [94] use temporal attention to determine the importance of each frame. For instance, [89] uses an attribute-driven method to disentangle features and re-weight frames, [72] selects the most discriminative frames automatically. These approaches, however, focus mainly on frame-level importance, neglecting the more fine-grained region-level importance. In [2] and [27], later frames are encouraged to focus on different regions from the previous frames, but it still requires substantial computational resources to extract deep features from each frame. Alternatively, we propose a pre-network patch selection mechanism that selects crucial patches from a video before processing them through a deep network, making the feature extraction process more efficient and effective. To ensure the effectiveness of our network, we propose warping the initial frame’s features to generate features for subsequent frames. While [28] generates pseudo features as well, its purposes and methods differ from ours significantly. [28] uses the vanilla

attention to generate the occlusion part features, while we propose a GOP-guided warping method to generate frame global features. Moreover, we present a self-adaptive strategy for achieving both low computation and high effectiveness, whereas the attention mechanism in [28] incurs higher computation costs. Another work related to our method is [85], which also uses a patch selection strategy in ReID to focus on the foreground area where the target person is located. However, it employs a completely different mechanism from ours. Specifically, in [85], deep network (ViT) features are first extracted for all patches in the entire image, and then key patches are selected based on these ViT features. As a result, each patch still needs to go through the full deep ViT for feature extraction, which incurs a significant computational cost. Our method addresses this issue by proposing a before-deep-network approach where key patches are selected prior to passing through the deep network for feature extraction, thus greatly reducing computational demands. Additionally, [85] only focuses on foreground-background separation, while our method further avoids redundant extraction of similar patches across different frames. With these novel designs, our method can be more efficient.

B. Efficient Vision Transformers and Token Pruning/Selection

In recent years, networks based on the transformer architecture have achieved tremendous success across various fields in computer vision [4], [102], [14], [74], [90], [25], [88], [95], [13], [99], [11], [32], [12]. However, the higher computational costs compared to traditional CNNs [24], [100], [96], [83], [98], [32], [52], [10], [8] have limited the application of these transformer methods in practical scenarios. To address this issue, researchers are actively seeking more efficient vision transformer approaches to reduce computational expenses. Some methods [81], [29], [23], [78], [79], [39] propose compact architectures to reduce the computational cost of self-attention. For example, [29] introduces an orthogonal self-attention mechanism that reduces the computational load while retaining more detailed information. Other approaches [82], [51], [15], [50] utilize quantization techniques to minimize computational overhead. For instance, [51] introduces a quantization approach specifically designed for vision transformers based on mixed-precision weights. Meanwhile, some methods [64], [70], [84], [42] employ distillation techniques to transfer the capabilities of large vision transformers to smaller transformers that require less computation.

Recently, the spatial redundancy in nature images has motivated some studies [62], [59], [57], [67], [54], [41], [34] to explore another kind of methods that discard nonessential tokens to improve processing speed. These techniques, known as token pruning or token selection, have been extended to temporal networks for efficient video processing [65], [38], [21], [17]. For example, [62] introduces a patch slimming method that discards useless patches in a top-down manner, with a novel mechanism that applies the crucial patches in a network's higher layers to guide the selection of patches in the shallower layers. [65] measures the importance of every region and chooses those with top scores to be used

for downstream processing. [40] proposes an automatic approach to prune unnecessary tokens that are semantically meaningless or distractive image backgrounds. [97] employs a detection module and scoring function to select only the class-discriminative regions for the computationally expensive texture feature extraction. [38] introduces a sparse video-text architecture with two forms of sparsity including edge sparsity to reduce the inter-token query-key communications and node sparsity to discard less-crucial tokens. [21] proposes a novel network that imitates human's sparse visual recognition in an end-to-end manner by representing images using only a very limited number of tokens in the latent space. Despite some success, these methods still exhibit two significant drawbacks. First, the pruning or selection mechanisms of these methods are not specifically designed for video-based person ReID, so directly applying them to this task may not guarantee task-optimal results. Second, most of these methods only focus on selecting sparse tokens, without further considering how to enhance the network's effectiveness with just sparse inputs. Unlike these methods, we propose a novel framework that resolves these issues, making the video ReID task more efficient and effective. Performance comparisons with these token pruning methods are presented in Sec. VI-B.

III. PRELIMINARIES

A. Group of Pictures

Group of Pictures (GOP) is a data structure encoded in compressed videos [68], [43], [53], [101], [19] to reduce the video storage cost. In a GOP, a video's first and subsequent frames are represented differently. The first frame, referred to as the **I-frame**, is represented by a fully encoded image $I \in \mathbb{R}^{3 \times H \times W}$. Each of the subsequent frames, known as the **P-frame**, is represented by a motion vector $M_t \in \mathbb{R}^{2 \times \frac{H}{16} \times \frac{W}{16}}$ combined with a residual map $R_t \in \mathbb{R}^{3 \times H \times W}$. The (i, j) -th pixel $M_t^{i,j}$ on M_t represents a coordinate, which indicates the displacement of the most similar patch in I-frame to the (i, j) -th patch in the t -th P-frame. $R_t^{:,16(i-1):16i,16(j-1):16j}$ represents residual errors by predicting the P-frame patch based on its I-frame displacement.

B. Patch-wise Representation

In transformer-based vision models like ViT, we generally represent the image patch-wise for sequential processing. Specifically, given an image with size of $\mathbb{R}^{3 \times H \times W}$, we first divide it equally into $H/16 \times W/16$ patches, each sized 16×16 , and then flatten it into the shape $\mathbb{R}^{N \times 768}$, where $N = H/16 \times W/16$, $768 = 3 \times 16 \times 16$ is the dimension of the flattened vector for each patch. To facilitate further operations, we reshape all maps within a GOP into a patch-based representation in this manner, getting I , M_t and R_t with the shape of $\mathbb{R}^{N \times 768}$, $\mathbb{R}^{2 \times N}$ and $\mathbb{R}^{N \times 768}$ respectively. Each $M_t^i \in \mathbb{R}^2$ in M_t represents a 2D spatial coordinate (h, w) , which we further convert into a single index v to correspond to the flattened coordinate, using the formula: $v = W/16 * h + w$. In this way, we get $M_t \in \mathbb{R}^N$.

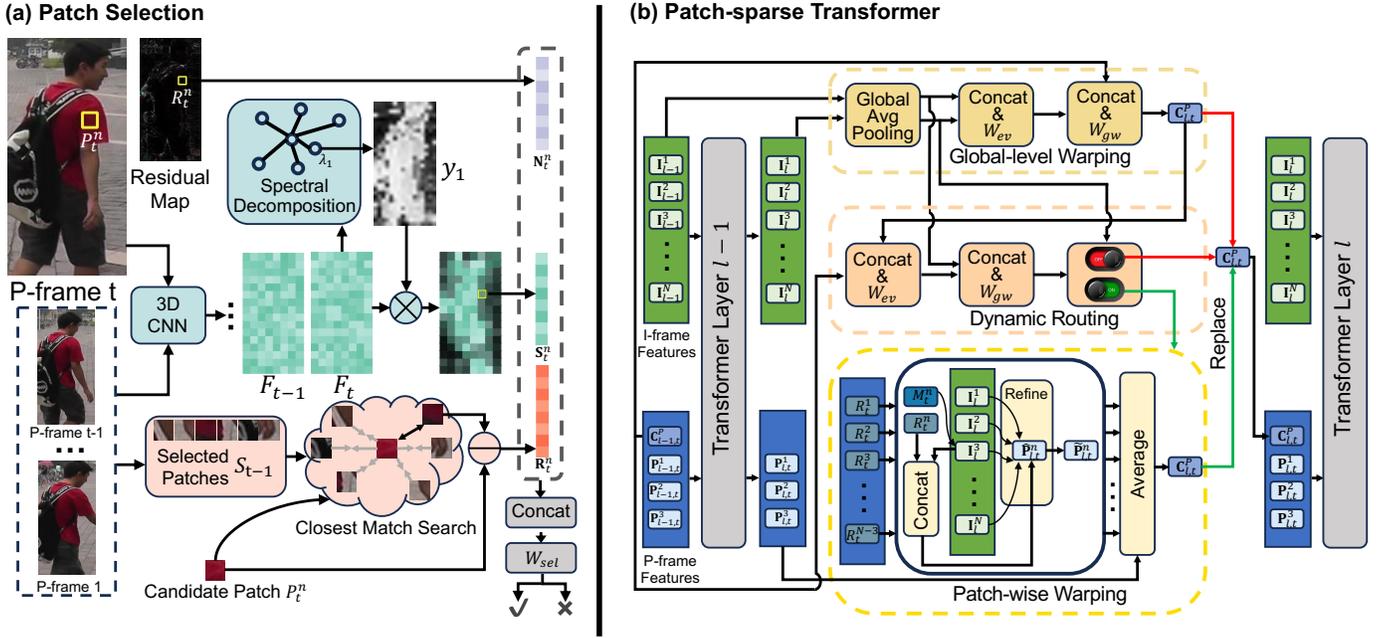


Fig. 1. **Overview of our framework**, including two components: (a) Patch Selection and (b) Patch-sparse Transformer. Note that, for simplicity of illustration, this figure only presents one I-frame with one P-frame in the patch-sparse transformer.

C. Spectral Decomposition

Spectral decomposition is a method proposed by [56] that can locate the prominent object in an image. Given an input feature $F \in \mathbb{R}^{H \times W}$, this method firsts computes an affinity matrix $A \in \mathbb{R}^{N \times N}$ ¹ through $A^{i,j} = F^i \cdot (F^j)^T$, where F^i refers to the i -th pixel on F , and then calculates the normalized Laplacian matrix L from A by using the formula $L = D^{-1/2}(D - A)D^{-1/2}$, where D is a diagonal matrix with $D^{i,i} = \sum_j A^{i,j}$. The eigenvectors of L are then computed and denoted as $\{y_0, y_1, \dots, y_{N-1}\}$, which are sorted according to their corresponding eigenvalues $\lambda_0 \leq \lambda_1 \leq \dots \leq \lambda_{N-1}$ in the ascending order. Based on the empirical observations in [56], the eigenvector $y_1 \in \mathbb{R}^N$ associated with the smallest nonzero eigenvalue λ_1 typically represents the most prominent object in a scene. In this work, we leverage this characteristic to develop a novel patch selection mechanism.

IV. PROPOSED METHOD

A. Overview

Following ViT [16], given a video clip $V \in \mathbb{R}^{3 \times H \times W \times T}$, our framework first transfers each frame to a patch-wise representation with size of $N \times 768$. Next, instead of extracting person features from all patches in all frames like previous transformer methods, we propose to use only a few patches selected from a video to be processed by the deep network. The model computation and feature redundancy can be thus greatly reduced. Specifically, we first propose a novel mechanism (Sec.IV-B) to automatically select patches. The selected patches are then fed into our proposed patch-sparse transformer (Sec.IV-C) to extract fine-grained person features,

which are ultimately used to identify persons. In the following sections, we describe each part of the framework in detail.

B. Patch Selection

The effectiveness of our framework is guaranteed by carefully selecting the best patches for feature extraction. Here, we utilize the GOP introduced in Sec.III that does not require additional sensors or computations to help this selection process. Specifically, we treat the I-frame within a GOP as the primary frame and use all patches within it as model inputs. We find retaining the complete I-frame is essential because it allows the extraction of important global features such as body shape, which is critical for the ReID task. Conversely, the remaining P-frames are selectively processed as subsidiary frames, from which only a subset of patches are selected for deep feature extraction.

To effectively and automatically select patches in the P-frames, we propose a new mechanism based on a scoring function as shown in Figure 1 (a) and Algorithm 1. Specifically, we denote the n -th patch in the t -th P-frame by P_t^n . For each P_t^n , we capture a selection-determinate feature for it, which is then inputted into a 3-layer MLP W_{sel} to generate a score. This score finally determines whether a patch should be selected or not. As the MLP input, the selection-determinate feature is the key to making this mechanism work. In the following subsection, we introduce this feature in detail.

1) *Selection-determinate Feature*: To guarantee an effective decision-making process with comprehensive considerations, we introduce two criteria as candidates for the selection-determinate feature: Patch Novelty N_t^n and Patch Semantics S_t^n . The motivation for using these two criteria is that we hope the selected patches are novel to the I-frame and not located in background areas, thus avoiding the redundant

¹ $N = H \times W$

feature extraction for the cross-frame repetitive and ReID-unrelated regions. Formally, \mathbf{N}_t^n is represented by the residual map R_t^n within a GOP, which measures the dissimilarity between the patch and the I-frame. \mathbf{S}_t^n represents a feature with semantic information related to the selection decision, which is obtained through a two-step process to extract and enhance features. Specifically, in the first step, we pass the input video through a shallow 4-layered 3D CNN to generate a feature map $F_t \in \mathbb{R}^{H/16 \times W/16 \times C}$ for each frame, which contains temporal and low-level information. In the second step, we utilize spectral decomposition to enhance F_t with saliency features, enabling it to capture foreground-background partitions more effectively. Specifically, we compute the eigenvector y_1 associated with the smallest nonzero eigenvalue λ_1 from F_t using the spectral decomposition method introduced in Sec.III-C. Through our empirical observations presented in Sec.VI-D2, y_1 typically captures the region of the target person in a ReID image. Therefore, we reshape y_1 to $\mathbb{R}^{H/16 \times W/16}$ and use it as a target-person indicator to enhance F_t by computing $\mathbf{S}_t = F_t \cdot y_1$. In this way, the resulting \mathbf{S}_t incorporates both patch semantics and foreground awareness, so it can be used to help W_{sel} to select better patches with crucial semantic information and within the target person.

2) *Progressive Selection*: Using W_{sel} in conjunction with input factors \mathbf{N} and \mathbf{S} , we can calculate patch scores for P-frames and select patches accordingly. The way of calculating scores for different P-frames is the following problem that needs to be carefully considered. Instead of calculating scores for all P-frames at once, we propose to use a progressive approach where patches are sequentially selected from the first frame to the last. This is because if we calculate scores in parallel, redundancy would occur when two similar patches originating from different P-frames are simultaneously chosen due to their closely matched input factors. Therefore, in our progressive approach, we construct a set \mathcal{S}_{t-1} containing all patches selected from the 1st to the current P-frames once the patch selection is completed in the $t - 1$ -th P-frame. When calculating the score for each patch P_t^n during the selection process of the subsequent frame t , we conduct a search for its closest match within \mathcal{S}_{t-1} ² and then calculate the pixel-wise residual \mathbf{R}_t^n between P_t^n and this matched patch. This residual information is incorporated as the third input to W_{sel} for scoring. By introducing \mathbf{R}_t^n into the scoring process, we make W_{sel} to have the awareness and adaptability to avoid selecting patches that closely resemble those already chosen, thereby effectively addressing the aforementioned redundancy issue by avoiding selecting similar patches repetitively.

3) *Differentiable Selection*: Eventually, we concatenate \mathbf{N}_t^n , \mathbf{S}_t^n and \mathbf{R}_t^n , which are inputted into W_{sel} to generate the score s_t^n . s_t^n is then used to make the selection decision. Although some past ReID works have explored methods to select key regions for emphasis, they either use a soft attention mechanism by assigning different importance weights to different patches and performing a weighted sum [66], which is differentiable but cannot achieve a hard decision, thus

unable to reduce computational load by excluding unnecessary patches; or they directly use a topK operation for patch selection [85], which can achieve a hard decision but is non-differentiable, therefore preventing the selection network from receiving gradients for updating. To overcome the limitations of these methods, we propose a novel mechanism to achieve differentiable and hard decisions by utilizing the saturating Sigmoid function [33] for patch selection. Specifically, during the training stage, we explore more decision space randomly by adding a standard Gaussian noise to s_t^n and getting \hat{s}_t^n . After that, we generate a binary score using $b_t^n = \mathbb{1}(\hat{s}_t^n > 0)$, and obtain a differentiable value d_t^n by:

$$d_t^n = \max(0, \min(1, 1.2\sigma(\hat{s}_t^n) - 0.1)), \quad (1)$$

where σ refers to the original Sigmoid function. Here, we use the binary score b_t^n for selection, i.e., selecting patch P_t^n if $b_t^n = 1$, and utilize the differentiable d_t^n to approximate gradients in back-propagation. Specifically, we use the gradient of d_t^n with respect to \hat{s}_t^n to serve as an approximation for the gradients required to update the parameters associated with the discrete gate b_t^n . The method for this gradient approximation can be implemented in PyTorch as follows: $b_t^n = b_t^n + d_t^n - d_t^n \cdot \text{detach}()$. Furthermore, during the inference phase, the process of sampling Gaussian noise is omitted. Instead, the discrete output is derived directly from the gate's original score and used for making decisions, i.e., selecting a patch if $s_t^n > 0$. In this way, the selection operation for P_t^n is completed, with a detailed description of the overall procedures shown in Algorithm 1. By using this method, we achieve effective selection that is both differentiable and capable of making hard decisions, thereby reducing computational load and enabling end-to-end training to simplify optimization.

C. Patch-sparse Transformer

1) *Motivation*: The previously mentioned method enables key patches to be selected from each P-frame. These patches can be input into a regular transformer like ViT for feature extraction with reduced computational cost. However, it is observed that when using the vanilla ViT as a feature extractor, employing these carefully selected patches as sparse inputs still results in a significant performance gap compared to using the entire video input. This gap could be attributed to the loss of global-aware information during the feature extraction process, which requires a non-local perspective so that cannot be captured when solely analyzing a limited number of patches. For example, without the complete view of a P-frame, it is difficult to identify whether a pattern in a selected patch is located on clothing or pants, which is a critical cue for person ReID. To address this issue, an intuitive and effective approach is to introduce global context so that the model can perceive global properties more effectively. However, trivial methods of obtaining context, such as global average pooling, are not directly applicable to our method since most P-frame positions do not have extracted features. As a result, context-based enhancement becomes a very challenging task in our framework. Fortunately, we observe that P-frames typically exhibit high similarities to the I-frame in a ReID video. This

²The patch similarity is measured using the same method employed in the H.264 video compression standard.

Algorithm 1 Patch Selection Algorithm

- 1: **Input:** video clip $V \in \mathbb{R}^{3 \times H \times W \times T}$ with T frames, including one I-frame $I \in \mathbb{R}^{3 \times H \times W}$ and $T - 1$ P-frames $\{P_t\}_{t=1}^{T-1} \in \mathbb{R}^{3 \times H \times W}$; residual map $R_t \in \mathbb{R}^{3 \times H \times W}$ in the GOP.
 - 2: Reshape I to $\mathbb{R}^{N \times 768}$, P_t to $\mathbb{R}^{N \times 768}$, R_t to $\mathbb{R}^{N \times 768}$.
 - 3: Pass the video V through a shallow 3D CNN, generating feature map $F_t \in \mathbb{R}^{H/16 \times W/16 \times C}$ for each P-frame P_t .
 - 4: Initialize \mathcal{S}_0 as an empty set.
 - 5: **for** t **in** $1, \dots, T - 1$ **do**
 - 6: **1. Obtain Patch Novelty** $\mathbf{N}_t = \{\mathbf{N}_t^n\}_{n=1}^N$:
 - 7: $\mathbf{N}_t^n \leftarrow R_t^n$.
 - 8: **2. Obtain Patch Semantics** $\mathbf{S}_t = \{\mathbf{S}_t^n\}_{n=1}^N$:
 - 9: Get $A \in \mathbb{R}^{N \times N}$ through $A^{i,j} = F_t^i \cdot (F_t^j)^T$ ($N = H/16 \times W/16$);
 - 10: Get $D \in \mathbb{R}^{N \times N}$ through $D^{i,j} = \sum_j A^{i,j}$;
 - 11: Get $L \in \mathbb{R}^{N \times N}$ through $L = D^{-1/2}(D - A)D^{-1/2}$;
 - 12: Compute the eigenvectors $\{y_0, y_1, \dots, y_{N-1}\}$ of L sorted according to their corresponding eigenvalues $\lambda_0 \leq \lambda_1 \dots \leq \lambda_{N-1}$ in the ascending order;
 - 13: Reshape y_1 to $\mathbb{R}^{H/16 \times W/16}$, compute $\mathbf{S}_t = F_t \cdot y_1$.
 - 14: **3. Obtain $\mathbf{R}_t = \{\mathbf{R}_t^n\}_{n=1}^N$ for progressive selection:**
 - 15: Search the most similar patch p in \mathcal{S}_{t-1} for P_t^n ;
 - 16: Compute pixel-wise residual \mathbf{R}_t^n between P_t^n and p .
 - 17: **4. Selection Decision for $\{P_t^n\}_{n=1}^N$:**
 - 18: Get s_t^n through $s_t^n = W_{sel}(\mathbf{N}_t^n || \mathbf{S}_t^n || \mathbf{R}_t^n)$;
 - 19: Select P_t^n if $s_t^n > 0$.
 - 20: **5. Update \mathcal{S} :**
 - 21: Get \mathcal{S}_t by appending all selected patches $\{P_t^n\}_{n=1}^{\hat{N}_t}$ from P_t to \mathcal{S}_{t-1} .
 - 22: **end for**
 - 23: **Return:** All patches $\{I^n\}_{n=1}^N$ in I-frame and the selected patches $\{\{P_t^n\}_{n=1}^{\hat{N}_t}\}_{t=1}^{T-1}$ from P-frames.
-

inspires us to propose a novel method that harnesses the I-frame features to generate pseudo global context $\mathbf{C}_{l,t}^P$ for each P-frame P_t in each transformer layer l and then use it to enhance global perception of the network. By doing so, as shown in Figure 1 (b), a novel patch-sparse transformer (PSFormer) can be developed to effectively handle sparse inputs (Sec.IV-C3). To simultaneously balance effectiveness and efficiency, we propose two complementary methods for generating pseudo context, namely patch-wise warping and global-level warping as follows:

(a) **Patch-wise Warping** generates pseudo context by first producing a pseudo feature for each unselected patch in P-frames. We notice that the GOP framework introduced in Sec.III can provide a perfect tool for this warping process. Specifically, in a GOP, the motion vector M represents the pairwise alignment between each patch of a P-frame and its displacement patch in the I-frame. According to the construction mechanism of GOP, these two patches are typically similar, with only minor differences being recorded in the residual map R . Taking advantage of this characteristic, we leverage the features of the displacement to generate a pseudo feature for the P-frame patch by warping using the residual

map that reflects the difference information. Formally, in the l -th layer of the transformer, we denote the feature for all I-frame patches as $\{\mathbf{I}_l^i\}_{i=1}^N$. And for each unselected patch P_t^n in the t -th P-frame, the transformer feature for its displacement patch is denoted as $\mathbf{I}_l^{M_t^n}$. As shown in the bottom part of Figure 1 (b), we combine $\mathbf{I}_l^{M_t^n}$ with the residual map R_t^n , and derive the pseudo feature $\hat{\mathbf{P}}_{l,t}^n$ for P_t^n through the operation $\hat{\mathbf{P}}_{l,t}^n = W_{pw}(\mathbf{I}_l^{M_t^n} || R_t^n)$ ³, where W_{pw} is a 3-layer MLP. In P-frames, the motion vector is predicted at the patch level, so it may not be the exact ground truth at the pixel-level. As a result, the alignment between P-frame patches and their displacements may not be perfect. To alleviate this potential error in alignment, we propose to assign different fusion importance weights to every location in $\mathbf{I}_l = \{\mathbf{I}_l^i\}_{i=1}^N$, thus refining $\hat{\mathbf{P}}_{l,t}^n$ by aggregating other potentially-related patches in the I-frame. Formally,

$$\tilde{\mathbf{P}}_{l,t}^n = \text{Softmax} \left(\frac{W_q(\hat{\mathbf{P}}_{l,t}^n)W_k(\mathbf{I}_l)^T}{\sqrt{d_k}} \right) W_v(\mathbf{I}_l), \quad (2)$$

where W_q , W_k and W_v are three fully-connected layers respectively aiming at producing the query, key and value features in cross-attention. Through this method, the weight assignment is implemented by measuring the correlation between $\hat{\mathbf{P}}_{l,t}^n$ and each $\mathbf{I}_l^i \in \{\mathbf{I}_l^i\}_{i=1}^N$, and feature aggregation is then completed by using these weights to perform a weighted summation over $\{W_v(\mathbf{I}_l^i)\}_{i=1}^N$. Finally, the frame global context $\mathbf{C}_{l,t}^P$ is calculated as the average over all patch features, including the transformer features $\mathbf{P}_{l,t}^n$ for selected patches and pseudo features $\hat{\mathbf{P}}_{l,t}^n$ for unselected patches.

(b) **Global-level Warping** directly warps the global context from the I-frame. As shown in the top part of Figure 1 (b), we first compute the average feature over all patches $\{\mathbf{I}_l^i\}_{i=1}^N$ in the I-frame, yielding a global context \mathbf{C}_l^I in the l -th layer of the transformer. Similarly, \mathbf{C}_{l-1}^I is obtained from the previous layer $l - 1$. We concatenate \mathbf{C}_l^I with \mathbf{C}_{l-1}^I , and obtain an evolution feature \mathbf{E}_l by:

$$\mathbf{E}_l = W_{ev}(\mathbf{C}_l^I || \mathbf{C}_{l-1}^I), \quad (3)$$

where W_{ev} is a 2-layer MLP. Through this process, \mathbf{E}_l can capture how the global context of the I-frame evolves from the previous to the current layer in the transformer. Due to the similarity between P-frames and I-frame, their global contexts display similar evolutionary patterns in a network. Therefore, we use \mathbf{E}_l as evolution guidance, generating the t -th P-frame's pseudo context $\mathbf{C}_{l,t}^P$ by warping $\mathbf{C}_{l-1,t}^P$ from the previous layer. Formally,

$$\mathbf{C}_{l,t}^P = W_{gw}(\mathbf{E}_l || \mathbf{C}_{l-1,t}^P), \quad (4)$$

where W_{gw} is a 2-layer MLP.

2) **Context Generation with Dynamic Routing:** The aforementioned two warping strategies have advantages and disadvantages. The patch-wise warping can provide better accuracy with its patch-level operation and the refinement mechanism, but it requires a higher computation cost. In contrast, the global-level warping is efficient, but may yield coarse and

³|| refers to concatenation.

inaccurate results as it overlooks spatial details and cross-frame differences. Moreover, when using global-level warping, each layer's $\mathbf{C}_{l,t}^P$ would be generated and thus inherit errors from the previous layer's $\mathbf{C}_{l-1,t}^P$. Consequently, the errors would accumulate in the transformer.

Based on the above discussions, we aim to strike a balance between the efficiency and effectiveness of these two methods, which is also the goal of our framework. We realize this target by proposing a novel dynamic routing mechanism to achieve self-adaptive feature generation, as shown in the middle part of Fig.1 (b). Specifically, in each transformer layer l , we first apply the global-level warping to generate a pseudo context $\mathbf{C}_{l,t}^P$, and then deploy an error-conditioned gate to determine whether a more refined version of $\mathbf{C}_{l,t}^P$ should be acquired through patch-wise warping. This gate is conditionally opened when the error amount accumulates to a certain extent. To measure the error amount, we employ global-level warping in a 'reverse' manner. To be specific, we use the pseudo context $\mathbf{C}_{l,t}^P, \mathbf{C}_{l-1,t}^P$ of the P-frame from two adjacent layers, generating an I-frame pseudo context $\hat{\mathbf{C}}_l^I$ for the l -th layer by warping from the context \mathbf{C}_{l-1}^I ⁴ of the previous layer $l-1$. Formally,

$$\hat{\mathbf{E}}_{l,t} = W_{ev}(\mathbf{C}_{l,t}^P || \mathbf{C}_{l-1,t}^P); \hat{\mathbf{C}}_l^I = W_{gw}(\hat{\mathbf{E}}_{l,t} || \mathbf{C}_{l-1}^I), \quad (5)$$

where W_{gw} and W_{ev} are the networks used in global-level warping (see Eq.3 and 4). Intuitively, increasing the level of errors accumulated in $\mathbf{C}_{l,t}^P$ results in a less effective evolution feature $\hat{\mathbf{E}}_{l,t}$, which leads to the coarser feature warping, making the warped $\hat{\mathbf{C}}_l^I$ from \mathbf{C}_{l-1}^I to exhibit a larger distance from the actual context \mathbf{C}_l^I . Therefore, we calculate the cosine distance $c_{l,t}$ between $\hat{\mathbf{C}}_l^I$ and \mathbf{C}_l^I , and use it as an approximate measurement of the error amount. $c_{l,t}$ is compared with a pre-defined activation threshold s . If $c_{l,t} > s$, we activate the gate, performing patch-level warping to obtain a refined pseudo context, which replaces $\mathbf{C}_{l,t}^P$. Otherwise, we keep the gate closed, directly using $\mathbf{C}_{l,t}^P$ from global-level warping for further operations. By implementing such a design, our network prioritizes using the coarse but computational-efficient global-level warping, while occasionally resorting to the more refined but computational-costly patch-wise warping only when the error accumulation reaches a certain threshold.

It is worth noting that the hyper-parameter s provides the flexibility to customize the model according to specific requirements. During the inference stage, if higher performance is the focus, the value of s can be reduced. On the other hand, if faster processing speed is desired, the value of s can be increased with faster processing but slightly sacrificed accuracy. Importantly, these adjustments can be made by only modifying the value of s , without the need to retrain the network. Experimental results of this dynamic adjustment are shown in Table VIII.

3) *Transformer Structure* : The pseudo context $\mathbf{C}_{l,t}^P$ obtained through dynamic routing is subsequently used to build our patch-spares transformer (PSFormer). Specifically,

⁴The I-frame context \mathbf{C}_l^I is computed as the average feature of all I-frame patches $\{\mathbf{I}_i^I\}_{i=1}^N$.

Algorithm 2 The l -th Layer of Patch-sparse Transformer

- 1: **Input**: output features from the $l-1$ -th layer, including $\{\mathbf{I}_i^I\}_{i=1}^N$ for all I-frame patches, $\{\{\mathbf{P}_{l,t}^n\}_{n=1}^{\hat{N}_t}\}_{t=1}^{T-1}$ for all selected P-frame patches; GOP; multi-head self attention MSA_l in the l -th layer.
 - 2: Process I-frame features: $\{\mathbf{I}_{l+1}^I\}_{i=1}^N = \text{MSA}_l(\{\mathbf{I}_i^I\}_{i=1}^N)$.
 - 3: **for** t **in** $1, \dots, T-1$ **do**
 - 4: **1. Global-level Warping**:
 - 5: $\mathbf{C}_l^I = \frac{1}{N} \sum_{i=1}^N \mathbf{I}_i^I$; $\mathbf{C}_{l-1}^I = \frac{1}{N} \sum_{i=1}^N \mathbf{I}_{i-1}^I$;
 - 6: $\mathbf{E}_l = W_{ev}(\mathbf{C}_l^I || \mathbf{C}_{l-1}^I)$; $\mathbf{C}_{l,t}^P = W_{gw}(\mathbf{E}_l || \mathbf{C}_{l-1,t}^P)$.
 - 7: **2. Dynamic Routing**:
 - 8: $\hat{\mathbf{E}}_{l,t} = W_{ev}(\mathbf{C}_{l,t}^P || \mathbf{C}_{l-1,t}^P)$; $\hat{\mathbf{C}}_l^I = W_{gw}(\hat{\mathbf{E}}_{l,t} || \mathbf{C}_{l-1}^I)$;
 - 9: $c_{l,t} = \text{Cos}(\hat{\mathbf{C}}_l^I, \mathbf{C}_l^I)$.
 - 10: **if** $c_{l,t} < s$ **then**
 - 11: $\{\mathbf{P}_{l+1,t}^n\}_{n=1}^{\hat{N}_t} = \text{MSA}_l(\{\mathbf{P}_{l,t}^n\}_{n=1}^{\hat{N}_t} || \mathbf{C}_{l,t}^P)$.
 - 12: **else**
 - 13: **3. Patch-wise Warping**:
 - 14: **for** n **in** $1, \dots, N - \hat{N}_t$ **do**
 - 15: Find the I-frame displacement $\mathbf{I}_l^{M_t^n}$ for the n -th unselected patch;
 - 16: $\hat{\mathbf{P}}_{l,t}^n = W_{pw}(\mathbf{I}_l^{M_t^n} || R_l^n)$;
 - 17: $\tilde{\mathbf{P}}_{l,t}^n = \text{Softmax}\left(\frac{W_q(\hat{\mathbf{P}}_{l,t}^n)W_k(\mathbf{I}_l)^T}{\sqrt{d_k}}\right)W_v(\mathbf{I}_l)$.
 - 18: **end for**
 - 19: Compute $\mathbf{C}_{l,t}^P$ as the average over all patch features (transformer features $\mathbf{P}_{l,t}^n$ for all selected patches and pseudo features $\hat{\mathbf{P}}_{l,t}^n$ for all unselected patches). In addition, average pooling all patch features to get 2×4 local tokens $\{\mathbf{L}_{l,t}^j\}_{j=1}^8$;
 - 20: $\{\mathbf{P}_{l+1,t}^n\}_{n=1}^{\hat{N}_t} = \text{MSA}_l(\{\mathbf{P}_{l,t}^n\}_{n=1}^{\hat{N}_t} || \mathbf{C}_{l,t}^P || \{\mathbf{L}_{l,t}^j\}_{j=1}^8)$.
 - 21: **end if**
 - 22: **end for**
 - 23: **Return**: $\{\mathbf{I}_{l+1}^I\}_{i=1}^N, \{\{\mathbf{P}_{l+1,t}^n\}_{n=1}^{\hat{N}_t}\}_{t=1}^{T-1}$
-

in the l -th layer of PSFormer, we denote the input features of all selected patches in the t -th P-frame as $\{\mathbf{P}_{l,t}^n\}_{n=1}^{\hat{N}_t}$. We append $\mathbf{C}_{l,t}^P$ as a token to $\{\mathbf{P}_{l,t}^j\}_{j=1}^{\hat{N}_t}$, and the resulted $[\mathbf{P}_{l,t}^1, \mathbf{P}_{l,t}^2, \dots, \mathbf{P}_{l,t}^{\hat{N}_t}, \mathbf{C}_{l,t}^P]$ is input into a multi-head self attention layer as in ViT. We observe that when patch-wise warping is used, the generated features not only encompass global context but also include the pseudo feature $\tilde{\mathbf{P}}_{l,t}^n$ for each unselected patch. To fully and efficiently utilize this fine-grained information, we employ average pooling to generate 2×4 tokens from these pseudo features. These tokens are also appended to $\{\mathbf{P}_{l,t}^n\}_{n=1}^{\hat{N}_t}$ as attention inputs. Consequently, PSFormer is constituted by successively stacking a total of L layers with the aforementioned structure. In Algorithm 2, we present the detailed procedures for each layer of our patch-sparse transformer.

After the last layer of the patch-sparse transformer, we employ a patch-wise warping process to generate pseudo features for all patches in the P-frames that have not been selected. Subsequently, features for all patches in all frames, including transformer output features for selected patches and pseudo features for unselected patches, are averaged to obtain

the final output of the patch-sparse transformer. This output feature is utilized for person re-identification. Note that our method can be applied to most of the existing video ReID models by replacing their used backbone network with our method that includes a patch selection stage followed by a patch-sparse transformer. These existing methods typically design modules to further capture more useful features, such as temporal information [66], [76], [2] like the walking pattern of pedestrians, based on the features extracted from the backbone. In such cases, the output of the patch-sparse transformer can be used as the backbone feature for further processing.

V. LOSS AND TRAINING

Following previous methods, we train our model using a combination of cross-entropy loss and hard triplet loss. In addition, we further propose an error-constraint loss function to ensure the effectiveness of our error measurement used in the dynamic routing mechanism. Specifically, in each transformer layer l , we randomly sample S numbers from the interval $[0, 1]$ and sort them ascendingly, yielding $\{\alpha_i\}_{i=1}^S$ with $\alpha_S > \alpha_{S-1} > \dots > \alpha_1$. Each α_i is used as a weight to generate a noisy I-frame global context $\tilde{C}_{l,i}^I = (1 - \alpha_i)C_l^I + \alpha_i\mathcal{N}$, where \mathcal{N} represents a standard Gaussian noise. Subsequently, using the global-level warping, we generate an evolution feature from $\tilde{C}_{l,i}^I$ and C_{l-1}^I , which is then used to obtain a reconstructed pseudo \tilde{C}_l^I by warping from C_{l-1}^I . We calculate the cosine distance $c_{l,i}$ between $\tilde{C}_{l,i}^I$ and the actual I-frame global context C_l^I . Based on the intuition that a higher α_i can lead to a noisier $\tilde{C}_{l,i}^I$ and finally result in a higher distance $c_{l,i}$, we formulate the following constraint:

$$\mathcal{L}_{error}^l = \sum_{i=1}^{S-1} \sum_{j=i+1}^S \max(0, c_{l,i} - c_{l,j}). \quad (6)$$

Using \mathcal{L}_{error} , we constrain the distance between the reconstructed I-frame features and the actual features to be positively correlated with the error amount. This constraint ensures that our error measurement method used in the dynamic routing mechanism can accurately reflect the accumulation condition of errors.

Using these losses, we adopt a 2-stage training process to optimize our method. In the first stage, we exclude the patch selection, feature warping, and dynamic routing components, using all patches from all frames as input to train the transformer layers in the patch-sparse transformer for 100 epochs. The loss function for this stage is the sum of cross-entropy loss \mathcal{L}_{cent} and hard triplet loss \mathcal{L}_{tri} , i.e., $\mathcal{L}_{cent} + \mathcal{L}_{tri}$. In the second stage, we further train the complete model for another 100 epochs. The loss function for this stage is the sum of cross-entropy loss \mathcal{L}_{cent} , hard triplet loss \mathcal{L}_{tri} and error-constraint loss \mathcal{L}_{error} , i.e., $\mathcal{L}_{cent} + \mathcal{L}_{tri} + \mathcal{L}_{error}$. Note that the training in the first stage is crucial. It allows the model to encounter more complete ReID images during training, which enables the model to learn better abilities to extract ReID-useful person features. We will show in experimental sections to demonstrate the effectiveness of our multi-stage training mechanism.

VI. EXPERIMENTS

A. Implementation Details

We use Adam as the optimizer with a momentum of 0.0005. The initial learning rate is set to 0.0005, which decays by 0.1 for every 40 epochs (overall 200 epochs). s in Sec. IV-C2 is set to 0.5. S in Eq. 6 is set to 4. We follow previous methods by using the restricted random sampling (RRS) strategy to sample 8 frames to generate each video clip, and the GOP that helps to perform the patch selection and patch-wise warping are obtained from this generated video clip after sampling. The frames are resized to 128×256 . We apply the commonly-used data augmentation strategies for training, including left-right flipping and random erasing. In PSFormer, the attention blocks follow the same setting as ViT with the same number of feature dimensions. We test and report the results of our method with 2 scales: Ours-small (Ours-S) with 8 layers and Ours-base (Ours-B) with 12 layers. We conduct experiments on the NVIDIA V100 GPUs.

B. Main Results

1) *Performance When Used as the Backbone for Existing Methods:* Table I presents the performance and computational costs when using our method as the backbone for existing methods across four datasets: MARS, LS-VID, iLiDS-VID, and PRID-2011. Specifically, we select four existing video-based ReID models, including MGH [75], SINet [2], GRL [49], and STMN[18], and replace their original backbone networks with the method proposed by us. We use two metrics to assess the model's performance: mean Average Precision (mAP) and rank-1, and two metrics to measure the computational cost of each method, including multiply-accumulate operations (MACs), which reflects the theoretical computation of the model; and milliseconds per video (ms/video), which calculates the total time for processing the whole 8 frames of each video on a realistic hardware platform (Tesla V100 in our experiments). The ms/video result is calculated as the average time for all test videos across all four datasets. Note that the reported computational cost of our method is the sum of the computational costs for both the patch selection stage and the patch-sparse transformer stage. In addition to our proposed method, we also test the method that uses two other widely-used networks as the model's backbone, including ResNet50 that is used by most of the previous methods, and the vanilla ViT without using our efficiency improvement strategies. For the comprehensive evaluation of method effectiveness, we validate and compare ViT and our method under two different sizes: ViT-B and Ours-B with 12 layers, and ViT-S and Ours-S with 8 layers. Based on the experimental results shown in Table I, we find that when used in conjunction with various existing methods, our proposed backbone can consistently achieve highly competitive results with only minimal computational requirements. Specifically, our approach (Ours-B) largely decreases computational cost by 74% compared to ViT-B and 28% compared to ResNet50, but can achieve results that closely resemble ViT-B and outperform ResNet50 significantly (+%5.1 mAP on mAP and +%8.5 on LS-VID). Additionally, with the same number of transformer

TABLE I
COMPARISON RESULTS WHEN USING FIVE DIFFERENT NETWORKS (RESNET50, ViT-B, OURS-B, ViT-S, OURS-S) AS THE BACKBONES FOR EXISTING VIDEO-BASED REID METHODS.

Method	Backbone	GMACs	ms/video	MARS		LS-VID		iLiDS-VID	PRID-2011
				mAP	rank-1	mAP	rank-1	rank-1	rank-1
ResNet50	ResNet50	32.7	94	81.0	88.1	68.5	80.4	83.3	90.2
ViT-B	ViT-B	88.9	272	86.7	90.0	78.3	87.4	90.0	93.8
Ours-B	Ours-B	23.5	78	86.1	89.5	77.7	86.9	89.3	93.4
ViT-S	ViT-S	51.1	144	85.4	89.2	77.0	85.8	88.5	93.1
Ours-S	Ours-S	14.4	45	84.8	89.1	76.5	85.4	88.0	92.8
MGH [75]	ResNet50	33.5	101	85.8	90.0	73.7	84.8	85.6	94.8
	ViT-B	89.7	280	88.0	91.0	79.5	88.8	91.6	96.2
	Ours-B	24.2	85	87.6	90.8	79.1	88.5	91.3	96.0
	ViT-S	52.0	153	87.1	90.6	78.7	88.0	91.0	95.7
	Ours-S	15.2	54	86.7	90.3	78.4	87.8	90.8	95.6
SINet [2]	ResNet50	37.7	115	86.3	91.2	79.6	87.6	92.7	96.5
	ViT-B	93.9	297	88.5	92.2	81.8	89.7	94.0	97.6
	Ours-B	28.6	100	88.2	92.1	81.7	89.5	93.7	97.6
	ViT-S	56.1	170	87.7	91.7	81.0	89.0	93.5	97.2
	Ours-S	19.4	70	87.5	91.6	80.8	88.8	93.3	97.1
GRL [49]	ResNet50	34.8	105	84.8	91.0	72.1	83.0	90.4	96.2
	ViT-B	91.0	285	87.9	91.9	79.1	88.5	92.9	97.2
	Ours-B	25.6	88	87.5	91.7	78.7	88.4	92.6	97.2
	ViT-S	53.3	159	87.0	91.5	78.0	87.7	92.0	97.0
	Ours-S	16.6	59	86.7	91.5	77.6	87.3	91.8	96.9
STMN[18]	ResNet50	36.4	110	84.5	90.5	69.5	82.1	91.5	95.5
	ViT-B	92.6	290	87.5	91.4	79.0	88.7	93.5	97.1
	Ours-B	27.2	94	87.3	91.4	78.6	88.4	93.4	97.0
	ViT-S	92.6	165	86.8	91.1	78.0	87.8	93.0	96.6
	Ours-S	17.2	64	86.6	90.9	77.7	87.5	92.9	96.4

TABLE II
COMPARISONS WITH OTHER EFFICIENT TRANSFORMER METHODS.

Method	GMACs	MARS		LS-VID	
		mAP	rank-1	mAP	rank-1
ViT-B [16]	88.9	86.7	90.0	78.3	87.4
DynamicViT [59]	35.5	83.8	88.4	74.8	84.7
Evit [41]	35.0	84.2	88.6	75.6	85.3
SPViT [34]	35.1	83.7	88.3	75.3	85.3
dTPS [67]	35.4	84.1	88.6	75.6	85.3
MViT [65]	31.9	84.8	89.0	76.6	85.9
DiffRate [7]	34.5	84.5	88.9	76.4	85.7
Ours-B	23.5	86.1	89.5	77.7	86.9

layers, Ours-S model also achieves results close to ViT-S while significantly reducing the computational cost. Moreover, compared to ViT-S, which has fewer layers, the Ours-B model, despite having more network layers, still requires less computational cost while achieving better performance. It is worth noting that due to the information loss associated with sparse input, our method exhibits a slight performance gap compared to ViT. However, given its only 26% computational requirement, our approach offers significantly higher practical utility value in real-time applications. Moreover, in comparison to the most-frequently-used ResNet50, our method not only reduces computational demands but also substantially enhances performance. Therefore, our method can be used to replace ResNet50 as a plug-and-play backbone, which can be seamlessly integrated with any other method to achieve both better performance and higher efficiency.

TABLE III
COMPARISON OF USING DIFFERENT EFFICIENT TRANSFORMER METHODS AS THE BACKBONES FOR EXISTING VIDEO-BASED REID METHODS.

Method	Backbone	GMACs	MARS		LS-VID	
			mAP	rank-1	mAP	rank-1
MGH [75]	DynamicViT [59]	36.2	85.8	89.7	76.5	85.5
	Evit [41]	35.7	86.2	90.1	77.0	86.2
	SPViT [34]	35.8	85.8	89.9	77.3	86.1
	dTPS [67]	36.1	86.4	90.1	77.4	86.4
	MViT [65]	32.6	86.7	90.3	78.1	86.9
	DiffRate [7]	34.7	86.4	90.1	77.9	86.8
	Ours-B	24.2	87.6	90.8	79.1	88.5
	Ours-B	24.2	87.6	90.8	79.1	88.5
GRL [49]	DynamicViT [59]	37.6	84.7	90.7	76.1	86.0
	Evit [41]	36.9	85.0	91.0	76.9	86.5
	SPViT [34]	37.0	85.0	90.9	77.2	86.4
	dTPS [67]	37.3	85.2	91.1	77.1	86.8
	MViT [65]	33.8	85.6	91.2	78.1	87.6
	DiffRate [7]	35.9	85.5	91.2	77.7	87.3
	Ours-B	25.6	87.5	91.7	78.7	88.4
	Ours-B	25.6	87.5	91.7	78.7	88.4

2) *Comparison with Other Efficient Transformers:* We further compare our approaches with other transformer methods that have the same objective of achieving high efficiency by using token pruning or selection techniques. To ensure fair comparisons, all these methods are evaluated using the same ViT-B network architecture and an equal number of layers. The results are presented in Table II, which includes both ReID performance and computational costs on the MARS and LS-VID datasets. The original ViT-B has the best performance, but it is very computationally intensive. Using token pruning or token selection techniques, the other methods can reduce computation while suffering from a slight loss in precision.

TABLE IV
COMPARISON WITH OTHER STATE-OF-THE-ART METHODS.

Method	Venue	MARS		iLiDS-VID	
		mAP	rank-1	rank-1	rank-5
GRL [49]	CVPR2021	84.8	91.0	90.4	98.3
CTL [46]	CVPR2021	86.7	91.4	89.7	97.0
STRF [1]	ICCV2021	86.1	90.3	89.3	-
STMN [18]	ICCV2021	84.5	90.5	91.5	98.5
PSTA [66]	ICCV2021	85.8	91.5	91.5	98.1
SINet [2]	CVPR2022	86.3	91.2	92.7	98.7
CAVIT [69]	ECCV2022	87.2	90.8	93.3	98.0
DCCT [48]	TNNLS2023	87.5	92.3	91.7	98.6
MSTAT [63]	TMM2023	85.3	91.8	93.3	99.3
FIDN [47]	TIP2023	86.8	91.5	91.3	98.0
Ours	-	88.2	92.1	93.7	99.5

This reduction in performance is primarily due to two reasons. First, the pruning or selection mechanisms of these methods are not specifically designed for video-based person ReID, so directly applying them to this task may not guarantee that the task-optimal patches are selected. Second, most of these methods focus only on selecting sparse tokens, without further considering how to address the loss of information resulting from pruning tokens. Compared to these methods, our method proposes to use ReID-task-tailored designs (patch novelty and patch semantics as selection-determinate features) for patch selection, so it can achieve the lowest computational costs while maintaining good performance by discarding redundant or unimportant patches to the greatest extent and preserving ReID-necessary patches as completely as possible. Additionally, we have innovatively designed a novel perception enhancement method for global information to address the issue of information loss when a lot of patches are discarded. By using these designs, we can still maintain good effectiveness while taking advantage of the reduced computational demand of sparse inputs. The results shown in Table II demonstrate that our method is both the most computationally efficient and performance-effective. Additionally, we also conduct performance comparisons when using these methods as the backbones for existing ReID methods. As shown in Table III, in this case, our method still achieves the best performance and highest efficiency. These results further demonstrate that our method is an excellent choice to serve as a plug-and-play backbone for other video ReID methods to achieve high efficiency and performance.

3) *Comparison with State-of-the-art Methods:* We further compare our method with other state-of-the-art video-based ReID methods to demonstrate our superior performance. All the methods used for comparison are advanced works published in top-tier conferences or journals in the past three years, and the reported results for ‘ours’ are obtained using SINet as the baseline with our method as its backbone. As shown in Table IV, across the four metrics on two datasets (MARS and iLiDS-VID), our method achieves the highest performance on three metrics (mAP for MARS, rank-1 and rank-5 for iLiDS-VID) and ranks second on the remaining one (rank-1 for MARS). These results highlight the high effectiveness and superiority of our method compared to the

TABLE V
ABLATION RESULTS OF PATCH SELECTION.

Method	mAP	rank-1	GMACs
Ours	86.1	89.5	23.5
Ours w/o patch selection	86.7	90.0	88.9
Ours w/o preserving whole I-frame	82.5	87.9	20.2
Ours w/o using patch novelty N	86.3	89.6	70.9
Ours w/o using patch semantics S	84.5	88.6	39.0
Ours w/o enhancement for S	85.6	89.0	31.5
Ours w/o frame-progressive selection (R)	86.1	89.5	33.4

previous SOTA methods. Notably, almost all the compared methods use conventional CNNs or transformer networks as the backbone for feature extraction, while our method significantly reduces the backbone’s computational demands through the innovatively proposed patch selection mechanism, which makes the network more efficient.

C. Ablation Study

We further perform ablation study to verify the effectiveness of our designs. Experiments in this section are conducted on MARS dataset based on the Ours-B model.

1) *Ablation of Patch Selection:* Table V presents the evaluation of our proposed patch selection mechanism. By removing patch selection and using all patches from a video for feature extraction, the computational cost is greatly increased by 278%. This result shows the significant role of our patch selection mechanism in reducing computation costs and achieving a high-efficient model. We further validate the contributions of various designs and components within our patch selection mechanism, including (1) the retention of the whole I-frame for feature extraction, (2) patch novelty **N** and patch semantics **S** as the selection-determinate features, (3) spectral decomposition methods used to enhance patch semantics **S**, and (4) the frame-progressive strategy with the factor **R** as a selection-determinate feature. The results indicate that removing any of these designs or components would result in either a significant decrease in model performance or a substantial increase in computational requirements. These results demonstrate that all designs in our mechanism can contribute to selecting patches that are more valuable and non-redundant, thus helping to achieve both high performance and reduced computational cost.

2) *Ablation of PSFormer:* Table VI presents an evaluation of different design choices in our PSFormer framework. (1) By excluding the use of context in each transformer layer, the mAP decreases by 3.2. This highlights the importance of incorporating global context during the feature extraction process, which helps to mitigate the problem of information loss caused by the sparse inputs of the network, thereby enhancing the representation ability of the features extracted by the transformer. (2) To generate the pseudo P-frame context feature, we employ a dynamic routing mechanism to select a feature warping method in each transformer layer. Without this mechanism and solely using global-level warping in all layers, the mAP decreases by 2.2. This is because global-level warping only utilizes a generalized global information for

TABLE VI
ABLATION RESULTS OF PSFORMER.

Method	mAP	rank-1	GMACs
Ours	86.1	89.5	23.5
Ours w/o using context in each layer	82.9	88.2	18.5
Ours solely using global-level warping	83.9	88.5	19.1
Ours solely using patch-wise warping	86.4	89.7	36.5

TABLE VII
DIFFERENT METHODS TO SELECT THE WARPING METHOD IN EACH NETWORK LAYER.

Method	mAP	rank-1	GMACs
Noise-conditioned Gate (Our method)	86.1	89.5	23.5
Randomly choosing in each layer	84.5	88.5	24.0
Using patch-wise warping every 4 layers	84.9	88.7	23.9

further processing, thus the obtained pseudo context lacks fine-grained detail features, which leads to the reduced accuracy. On the other hand, solely using patch-wise warping achieves slightly better performance than our method by utilizing more fine-grained information, but it requires a significantly higher computation cost (+71%) due to the patch-level operations and feature fusion mechanism. As demonstrated by the experimental results, by dynamically and adaptively selecting between two complementary warping methods using the proposed dynamic routing mechanism, our method can achieve the optimal balance between low computation and high effectiveness.

3) *Effectiveness of Noise-conditioned Gate*: The proposed dynamic routing mechanism in our method employs a noise-conditioned gate to select between the global-level warping and patch-wise warping in each network layer. To verify the effectiveness of this gate, we compare it with two other strategies and present the results in Table VII. Specifically, the first comparison strategy uses a random warping method in each layer instead of using the proposed selection method, which decreases the mAP by 1.6. The second comparison strategy is to use patch-wise warping every four layers. This method also decreases the mAP by 1.2. These results demonstrate that, compared to random or manually designed methods, the proposed noise-conditioned gate can help the network to select the more appropriate warping method in each layer automatically and self-adaptively, thus achieving a better balance between computation and performance.

4) *Dynamic Adjustment for Computation-Performance Balance*: As introduced detailedly in Sec. IV-C2, in the proposed dynamic routing mechanism, we introduce a hyper-parameter s to control the activation of the noise-conditioned gate. We find that by setting s to different values, we can dynamically adjust the balance between model computation and ReID performance based on user requirements. Specifically, if higher performance is the user’s focus, the value of s can be decreased, which results in more network layers using the more accurate but computationally demanding patch-wise warping, thus leading to better performance but more computation cost. Conversely, if a faster processing speed is desired, the value

TABLE VIII
PERFORMANCE AND COMPUTATION FOR DIFFERENT THRESHOLDS s .

Threshold s	0.4	0.5	0.6	0.7	0.8	0.9
mAP	86.3	86.1	85.5	85.0	84.6	84.2
rank-1	89.5	89.5	89.1	88.9	88.7	88.5
GMACs	27.1	23.5	21.9	20.9	19.9	19.5

TABLE IX
ABLATION RESULTS OF LOSS FUNCTIONS AND TRAINING METHODS.

Method	mAP	rank-1
Ours	86.1	89.5
Ours w/o \mathcal{L}_{cent}	79.2	85.8
Ours w/o \mathcal{L}_{tri}	82.9	87.9
Ours w/o \mathcal{L}_{error}	84.9	88.5
Ours w/o training in the first stage	83.3	88.0

of s can be increased to achieve faster processing but with a slight sacrifice in accuracy by using global-level warping more frequently. In Table VIII, we show how the model computation and performance change as s increases, which demonstrates the high flexibility and controllability of our approach that is valuable in real-world applications.

5) *Evaluation of Loss and Training Methods*: We further evaluate the effectiveness of the proposed loss functions and two-stage training methods. The results are presented in Table IX. As introduced in Sec. V, the loss function used to optimize our method is the sum of three components: cross-entropy loss \mathcal{L}_{cent} , hard triplet loss \mathcal{L}_{tri} , and error-constraint loss \mathcal{L}_{error} . As shown in Table IX, removing any one of these three components will significantly reduce model performance, which demonstrates the rationality of using all three functions together. Notably, the error-constraint loss is innovatively proposed by us and is specifically designed for the dynamic routing mechanism in our PSFormer. It ensures that the error measurement method used in the mechanism can accurately reflect the error accumulation condition. The drop in performance when removing \mathcal{L}_{error} demonstrates the importance of this newly proposed loss function. Furthermore, we validate the two-stage training method used for optimizing our model. As discussed in Sec V, in the first stage, we use all patches from all frames as input to only train the transformer layers in the PSFormer. In the second stage, we further train the complete model with the patch selection modules. The training in the first stage is crucial since it allows the model to encounter more complete ReID images, which enables the model to learn better abilities to extract ReID-useful person features. Thus, compared to the training method only with the second stage, our two-stage training method can increase the mAP by 2.8.⁵ In addition, we present the training curves to further demonstrate the advantages of our 2-stage training method. As shown in Figure 3, compared to training for 200 epochs using only the second stage, training with the first stage for 100 epochs followed by the second stage for another

⁵For fair comparisons, we use the same total number of epochs for training the 1-stage method and 2-stage method.



Fig. 2. Visualization of selected patches for the 1st, 3rd, 5th, and 7th frames in a video clip. The selected patches are marked by yellow masks. (Best viewed in color)

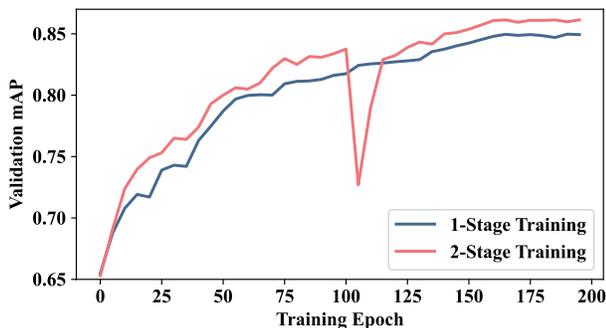


Fig. 3. Convergence curves on the MARS dataset.

100 epochs in our method achieves better convergence. Note that the training curve of our 2-stage method shows a sudden change at the 100th epoch, as this is when the second stage begins, and additional modules are introduced (including the patch selection modules, as well as the feature warping and dynamic routing components in the patch-sparse transformer). However, the model quickly converges afterward, ultimately achieving better performance than the single-stage method. These results demonstrate the rationality and significant effectiveness of our proposed training methods.

D. Further Analysis

1) *Visualizations of Selected Patches:* To gain a deeper and more intuitive understanding of our patch selection method, in Figure 2, we present some examples of the selected patches across different frames of a video clip. Specifically, we sample the 1st, 3rd, 5th, and 7th frames in a video clip for presentation, and the selected patches in these frames are highlighted in yellow. A video clip’s first frame is regarded as the I-frame, and in our method, all patches within it are selected to extract global information. While for the following P-frames,

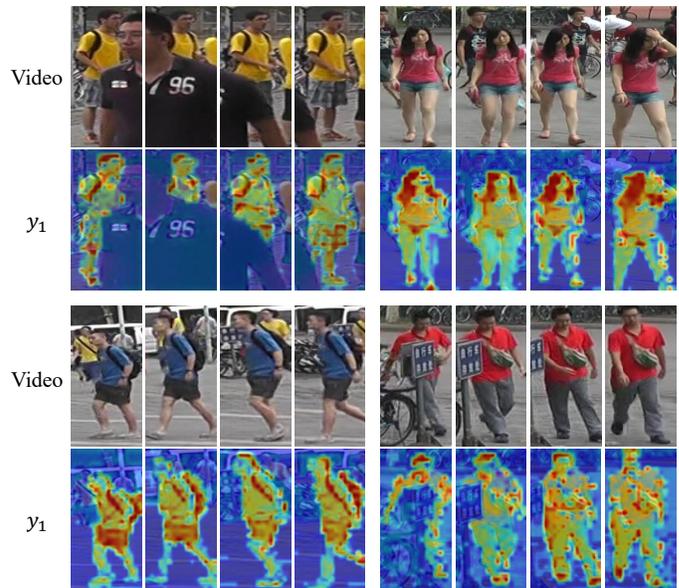


Fig. 4. Visualization of the eigenvector y_1 associated with the smallest nonzero eigenvalue obtained through the spectral decomposition method illustrated in Sec.III-C

we observe that only a few key patches are chosen while most redundant patches repeating from previous frames or located in background regions are effectively pruned. Interestingly, we find that a small number of patches that repeat across multiple frames are still consistently reselected. A reason behind this phenomenon could be that these patches are typically situated in crucial areas for ReID such as the human face, so the repetitive feature extraction for them can effectively emphasize and enhance the ReID-important information. Furthermore, as shown in the last row of Figure 2, we observe that in a challenging condition when the target person is occluded by other persons in later P-frames of a video clip, our proposed method is still capable of making the correct patch selection decisions by excluding and not choosing the regions where these occluded persons are located. This is because our patch selection method uses patch semantics information as a selection-determinant feature, which is obtained by passing the video clip through a shallow 3D-CNN network. As such, it contains temporal information that can help differentiate the target person from the occluded persons in the video. These results demonstrate that our method can effectively select the most appropriate patches even in complex videos, thereby greatly reducing the computational cost of the subsequent transformer to realize a more efficient video ReID backbone.

2) *Effectiveness of Spectral Decomposition:* As detailed in Sec.IV-B1, to select only foreground target person regions for deep feature extraction, we employ spectral decomposition to compute an eigenvector associated with the smallest nonzero eigenvalue. This eigenvector can represent the most prominent object in a scene, which, in a ReID image, is typically the target person. To demonstrate this, we further conduct the following visual verification and quantitative validations: (1) **Visual verification.** As shown in Figure 4, we visualize the eigenvectors y_1 associated with the smallest nonzero eigen-

TABLE X
ABLATION STUDY FOR TRAINING EPOCHS IN DIFFERENT STAGES.

Stage1 Epoch	Stage2 Epoch	MARS mAP	MARS rank-1
70	130	86.00	89.34
80	120	86.12	89.50
90	110	86.18	89.50
100	100	86.13	89.53
110	90	86.08	89.46
120	80	86.11	89.43
130	70	86.03	89.39

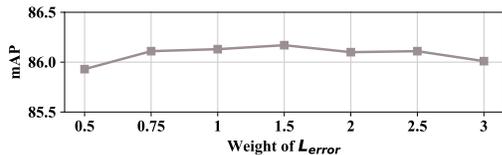


Fig. 5. Ablation study for the weight of error-constraint loss \mathcal{L}_{error} on MARS dataset.

value from several ReID video frames. Specifically, a heat map is generated based on the value of y_1 in each patch for visualization. It can be seen that these eigenvectors effectively capture the region of the target person, even in frames where the background contains distracting people or objects. This benefits from the temporal information contained in F_t , which serves as the input feature for spectral decomposition and can help accurately distinguish the target person from background objects. These visualization results demonstrate the effectiveness of our method. (2) **Quantitative validation.** With the assistance of the Segment Anything Model [60], we obtain pixel-level ground truth masks for the region of the target person in each video frame from the MARS validation set. We then apply our method to compute the eigenvector associated with the smallest nonzero eigenvalue for each frame, generate a binary mask using a threshold of 0, and calculate its mIoU with the corresponding ground truth mask. Across the entire dataset, our method achieves a high mIoU of 74.1%. This result reveals that the prominent object represented by these eigenvectors overlaps highly with the target person, further demonstrating the high effectiveness of our method.

3) **Hyperparameter Analysis:** We further validate the robustness of our method to two key hyperparameter settings: (1) the number of training epochs in each stage, and (2) the weight of each loss function. To be specific, we keep the total number of training epochs at 200, vary the number of training epochs in the first stage to range from 70 to 130, and adjust the number of epochs in the second stage accordingly. As shown by the results presented in Table X, our method consistently achieves stable and excellent performance across all evaluated settings on the MARS dataset. We then evaluate the performance of our method under varying weights for each loss function. Specifically, in the second training stage, our approach employs a weighted sum of three loss functions for optimization: cross-entropy loss \mathcal{L}_{cent} , hard triplet loss \mathcal{L}_{tri} , and error-constraint loss \mathcal{L}_{error} . For \mathcal{L}_{cent} and \mathcal{L}_{tri} , we directly follow most previous works [2], [75], [66] by setting their weights to 1. Therefore, we primarily focus on the weight W_{error} for our newly proposed function \mathcal{L}_{error} . The results

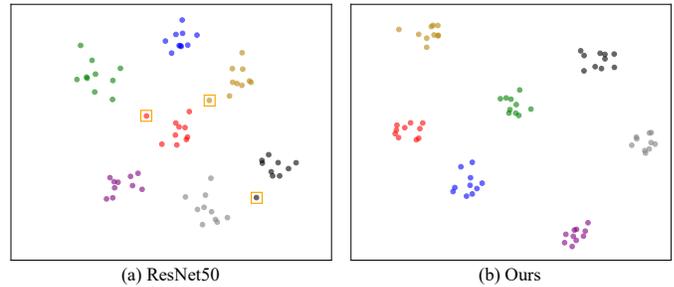


Fig. 6. t-SNE visualization results of features from ResNet50 and our method. Circles marked by the orange bounding boxes refer to videos likely to be matched to the incorrect identities.

presented in Figure 5 indicate that our method maintains stable performance across a wide range of W_{error} values from 0.5 to 3. These results suggest that our method is not sensitive to the hyperparameter settings, since it can consistently achieve stable and excellent performance across a broad spectrum of hyperparameter values, demonstrating its high effectiveness and significant robustness.

4) **t-SNE Visualization Analysis:** To further validate the effectiveness of our method, we provide the t-SNE visualizations of network features in Figure 6. In this visualization, each color represents a different identity, and each circle corresponds to features extracted from a video clip. It is observed that with the most-commonly-used ResNet50 as the backbone, features of different videos from the same identity exhibit significant intra-identity variation and are not tightly clustered. Moreover, the distances between distributions of different identities are relatively close, which may cause some videos at the distribution boundaries—such as the ones marked by orange bounding boxes—to be matched to the incorrect identities. In contrast, our method shows clear advantages: it significantly reduces intra-identity variation and enhances the distinctness of inter-identity boundaries. This improvement contributes to more accurate person re-identification, demonstrating the superior effectiveness of our approach compared to ResNet50, even with lower computational costs.

VII. CONCLUSION

This paper proposes a new effective and efficient plug-and-play backbone for video person ReID. This backbone consists of two components: a patch selection mechanism to prune redundant patches for reduced computational load, and a patch-sparse transformer that achieves high-performance feature extraction enhanced by pseudo frame-global context. In comparison to the widely used ResNet50 backbone, our approach significantly reduces computational costs while effectively improving performance. We regard our method as a general and practical approach that can be utilized in real-world applications.

REFERENCES

- [1] Abhishek Aich, Meng Zheng, Srikrishna Karanam, Terrence Chen, Amit K Roy-Chowdhury, and Ziyang Wu. Spatio-temporal representation factorization for video-based person re-identification. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 152–162, 2021.

- [2] Shutao Bai, Bingpeng Ma, Hong Chang, Rui Huang, and Xilin Chen. Salient-to-broad transition for video person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7339–7348, 2022.
- [3] Chengzhi Cao, Xuexiang Fu, Hongjian Liu, Yukun Huang, Kunyu Wang, Jiebo Luo, and Zheng-Jun Jia. Event-guided person re-identification via sparse-dense complementary learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17990–17999, 2023.
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020.
- [5] Dapeng Chen, Hongsheng Li, Tong Xiao, Shuai Yi, and Xiaogang Wang. Video person re-identification with competitive snippet-similarity aggregation and co-attentive snippet embedding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1169–1178, 2018.
- [6] Guangyi Chen, Yongming Rao, Jiwen Lu, and Jie Zhou. Temporal coherence or temporal motion: Which is more critical for video-based person re-identification? In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VIII 16*, pages 660–676. Springer, 2020.
- [7] Mengzhao Chen, Wenqi Shao, Peng Xu, Mingbao Lin, Kaipeng Zhang, Fei Chao, Rongrong Ji, Yu Qiao, and Ping Luo. DiffRate: Differentiable compression rate for efficient vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, 2023.
- [8] Tianrun Chen, Chaotao Ding, Lanyun Zhu, Ying Zang, Yiyi Liao, Zejian Li, and Lingyun Sun. Reality3dsketch: Rapid 3d modeling of objects from single freehand sketches. *IEEE Transactions on Multimedia*, 2023.
- [9] Tianlong Chen, Shaojin Ding, Jingyi Xie, Ye Yuan, Wuyang Chen, Yang Yang, Zhou Ren, and Zhangyang Wang. Abd-net: Attentive but diverse person re-identification. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8351–8361, 2019.
- [10] Tianrun Chen, Chenglong Fu, Lanyun Zhu, Papa Mao, Jia Zhang, Ying Zang, and Lingyun Sun. Deep3dsketch: 3d modeling from free-hand sketches with view-and structural-aware adversarial training. *arXiv preprint arXiv:2312.04435*, 2023.
- [11] Tianrun Chen, Ankang Lu, Lanyun Zhu, Chaotao Ding, Chunan Yu, Deyi Ji, Zejian Li, Lingyun Sun, Papa Mao, and Ying Zang. Sam2-adapter: Evaluating & adapting segment anything 2 in downstream tasks: Camouflage, shadow, medical image segmentation, and more. *arXiv preprint arXiv:2408.04579*, 2024.
- [12] Tianrun Chen, Chunan Yu, Jing Li, Jianqi Zhang, Lanyun Zhu, Deyi Ji, Yong Zhang, Ying Zang, Zejian Li, and Lingyun Sun. Reasoning3d-grounding and reasoning in 3d: Fine-grained zero-shot open-vocabulary 3d reasoning part segmentation via large vision-language models. *arXiv preprint arXiv:2405.19326*, 2024.
- [13] Tianrun Chen, Lanyun Zhu, Chaotao Ding, Runlong Cao, Yan Wang, Shangzhan Zhang, Zejian Li, Lingyun Sun, Ying Zang, and Papa Mao. Sam-adapter: Adapting segment anything in underperformed scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3367–3375, 2023.
- [14] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1290–1299, 2022.
- [15] Yifu Ding, Haotong Qin, Qinghua Yan, Zhenhua Chai, Junjie Liu, Xiaolin Wei, and Xianglong Liu. Towards accurate post-training quantization for vision transformer. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 5380–5388, 2022.
- [16] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [17] Matthew Dautson, Yin Li, and Mohit Gupta. Eventful transformers: Leveraging temporal redundancy in vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, 2023.
- [18] Chanho Eom, Geon Lee, Junghyup Lee, and Bumsu Ham. Video-based person re-identification with spatial and temporal memory networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12036–12045, 2021.
- [19] Zhipeng Fan, Jun Liu, and Yao Wang. Motion adaptive pose estimation from compressed videos. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11719–11728, 2021.
- [20] Yang Fu, Xiaoyang Wang, Yunchao Wei, and Thomas Huang. Sta: Spatial-temporal attention for large-scale video-based person re-identification. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 8287–8294, 2019.
- [21] Ziteng Gao, Zhan Tong, Limin Wang, and Mike Zheng Shou. Sparseformer: Sparse visual recognition via limited latent tokens. *arXiv preprint arXiv:2304.03768*, 2023.
- [22] Xinqian Gu, Hong Chang, Bingpeng Ma, Hongkai Zhang, and Xilin Chen. Appearance-preserving 3d convolution for video-based person re-identification. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 228–243. Springer, 2020.
- [23] Dongchen Han, Xuran Pan, Yizeng Han, Shiji Song, and Gao Huang. Flatten transformer: Vision transformer using focused linear attention. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5961–5971, 2023.
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [25] Shuting He, Hao Luo, Pichao Wang, Fan Wang, Hao Li, and Wei Jiang. Transreid: Transformer-based object re-identification. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 15013–15022, 2021.
- [26] Ruibing Hou, Hong Chang, Bingpeng Ma, Rui Huang, and Shiguang Shan. Bicnet-tks: Learning efficient spatial-temporal representation for video person re-identification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2014–2023, 2021.
- [27] Ruibing Hou, Hong Chang, Bingpeng Ma, Shiguang Shan, and Xilin Chen. Temporal complementary learning for video person re-identification. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXV 16*, pages 388–405. Springer, 2020.
- [28] Ruibing Hou, Bingpeng Ma, Hong Chang, Xinqian Gu, Shiguang Shan, and Xilin Chen. Vrstc: Occlusion-free video person re-identification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7183–7192, 2019.
- [29] Huaibo Huang, Xiaoqiang Zhou, and Ran He. Orthogonal transformer: An efficient vision transformer backbone with token orthogonalization. *Advances in Neural Information Processing Systems*, 35:14596–14607, 2022.
- [30] Meiyuan Huang, Chunging Hou, Qingyuan Yang, and Zhipeng Wang. Reasoning and tuning: Graph attention network for occluded person re-identification. *IEEE Transactions on Image Processing*, 32:1568–1582, 2023.
- [31] Yan Huang, Zhang Zhang, Qiang Wu, Yi Zhong, and Liang Wang. Enhancing person re-identification performance through in vivo learning. *IEEE Transactions on Image Processing*, 2023.
- [32] Deyi Ji, Feng Zhao, Lanyun Zhu, Wenwei Jin, Hongtao Lu, and Jieping Ye. Discrete latent perspective learning for segmentation and detection. *arXiv preprint arXiv:2406.10475*, 2024.
- [33] Lukasz Kaiser, Samy Bengio, Aurko Roy, Ashish Vaswani, Niki Parmar, Jakob Uszkoreit, and Noam Shazeer. Fast decoding in sequence models using discrete latent variables. In *International Conference on Machine Learning*, pages 2390–2399. PMLR, 2018.
- [34] Zhenglun Kong, Peiyan Dong, Xiaolong Ma, Xin Meng, Wei Niu, Mengshu Sun, Xuan Shen, Geng Yuan, Bin Ren, Hao Tang, et al. Spvit: Enabling faster vision transformers via latency-aware soft token pruning. In *European Conference on Computer Vision*, pages 620–640. Springer, 2022.
- [35] Jianing Li, Jingdong Wang, Qi Tian, Wen Gao, and Shiliang Zhang. Global-local temporal representations for video person re-identification. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3958–3967, 2019.
- [36] Jianing Li, Shiliang Zhang, and Tiejun Huang. Multi-scale 3d convolution network for video based person re-identification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8618–8625, 2019.
- [37] Yulin Li, Jianfeng He, Tianzhu Zhang, Xiang Liu, Yongdong Zhang, and Feng Wu. Diverse part discovery: Occluded person re-identification with part-aware transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2898–2907, 2021.
- [38] Yi Li, Kyle Min, Subarna Tripathi, and Nuno Vasconcelos. Svtt: Temporal learning of sparse video-text transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18919–18929, 2023.
- [39] Yanyu Li, Geng Yuan, Yang Wen, Ju Hu, Georgios Evangelidis, Sergey Tulyakov, Yanzhi Wang, and Jian Ren. Efficientformer: Vision transformers at mobilenet speed. *Advances in Neural Information*

- Processing Systems*, 35:12934–12949, 2022.
- [40] Youwei Liang, GE Chongjian, Zhan Tong, Yibing Song, Jue Wang, and Pengtao Xie. Evit: Expediting vision transformers via token reorganizations. In *International Conference on Learning Representations*, 2021.
- [41] Youwei Liang, Chongjian Ge, Zhan Tong, Yibing Song, Jue Wang, and Pengtao Xie. Not all patches are what you need: Expediting vision transformers via token reorganizations. In *International Conference on Learning Representations*, 2022.
- [42] Han Lin, Guangxing Han, Jiawei Ma, Shiyuan Huang, Xudong Lin, and Shih-Fu Chang. Supervised masked knowledge distillation for few-shot transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19649–19659, 2023.
- [43] Ting-Lan Lin, Sandeep Kanumuri, Yuan Zhi, David Poole, Pamela C Cosman, and Amy R Reibman. A versatile model for packet loss visibility and its application to packet prioritization. *IEEE Transactions on Image Processing*, 19(3):722–735, 2009.
- [44] Xinyu Lin, Jinxing Li, Zeyu Ma, Huafeng Li, Shuang Li, Kaixiong Xu, Guangming Lu, and David Zhang. Learning modal-invariant and temporal-memory for video-based visible-infrared person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20973–20982, 2022.
- [45] Hao Liu, Zequn Jie, Karlekar Jayashree, Meibin Qi, Jianguo Jiang, Shuicheng Yan, and Jiashi Feng. Video-based person re-identification with accumulative motion context. *IEEE transactions on circuits and systems for video technology*, 28(10):2788–2802, 2017.
- [46] Jiawei Liu, Zheng-Jun Zha, Wei Wu, Kecheng Zheng, and Qibin Sun. Spatial-temporal correlation and topology learning for person re-identification in videos. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4370–4379, 2021.
- [47] Liangchen Liu, Xi Yang, Nannan Wang, and Xinbo Gao. Frequency information disentanglement network for video-based person re-identification. *IEEE Transactions on Image Processing*, 2023.
- [48] Xuehu Liu, Chenyang Yu, Pingping Zhang, and Huchuan Lu. Deeply coupled convolution–transformer with spatial–temporal complementary learning for video-based person re-identification. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [49] Xuehu Liu, Pingping Zhang, Chenyang Yu, Huchuan Lu, and Xiaoyun Yang. Watching you: Global-guided reciprocal learning for video-based person re-identification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13334–13343, 2021.
- [50] Yijiang Liu, Huanrui Yang, Zhen Dong, Kurt Keutzer, Li Du, and Shanghang Zhang. Noisyquant: Noisy bias-enhanced post-training activation quantization for vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20321–20330, 2023.
- [51] Zhenhua Liu, Yunhe Wang, Kai Han, Wei Zhang, Siwei Ma, and Wen Gao. Post-training quantization for vision transformer. *Advances in Neural Information Processing Systems*, 34:28092–28103, 2021.
- [52] Zhikang Liu and Lanyun Zhu. Label-guided attention distillation for lane segmentation. *Neurocomputing*, 438:312–322, 2021.
- [53] Guo Lu, Tianxiong Zhong, Jing Geng, Qiang Hu, and Dong Xu. Learning based multi-modality image and video compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6083–6092, 2022.
- [54] Haoyu Ma, Zhe Wang, Yifei Chen, Deying Kong, Liangjian Chen, Xingwei Liu, Xiangyi Yan, Hao Tang, and Xiaohui Xie. Ppt: token-pruned pose transformer for monocular and multi-view human pose estimation. In *European Conference on Computer Vision*, pages 424–442. Springer, 2022.
- [55] Niall McLaughlin, Jesus Martinez Del Rincon, and Paul Miller. Recurrent convolutional network for video-based person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1325–1334, 2016.
- [56] Luke Melas-Kyriazi, Christian Rupprecht, Iro Laina, and Andrea Vedaldi. Deep spectral methods: A surprisingly strong baseline for unsupervised semantic segmentation and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8364–8375, 2022.
- [57] Lingchen Meng, Hengduo Li, Bor-Chun Chen, Shiyi Lan, Zuxuan Wu, Yu-Gang Jiang, and Ser-Nam Lim. Adavit: Adaptive vision transformers for efficient image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12309–12318, 2022.
- [58] Yongming Rao, Guangyi Chen, Jiwen Lu, and Jie Zhou. Counterfactual attention learning for fine-grained visual categorization and re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1025–1034, 2021.
- [59] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit: Efficient vision transformers with dynamic token sparsification. *Advances in neural information processing systems*, 34:13937–13949, 2021.
- [60] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024.
- [61] Arulkumar Subramaniam, Athira Nambiar, and Anurag Mittal. Co-segmentation inspired attention networks for video-based person re-identification. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 562–572, 2019.
- [62] Yehui Tang, Kai Han, Yunhe Wang, Chang Xu, Jianyuan Guo, Chao Xu, and Dacheng Tao. Patch slimming for efficient vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12165–12174, 2022.
- [63] Ziyi Tang, Ruimao Zhang, Zhanglin Peng, Jinrui Chen, and Liang Lin. Multi-stage spatio-temporal aggregation transformer for video person re-identification. *IEEE Transactions on Multimedia*, 25:7917–7929, 2022.
- [64] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347–10357. PMLR, 2021.
- [65] Junke Wang, Xitong Yang, Hengduo Li, Li Liu, Zuxuan Wu, and Yu-Gang Jiang. Efficient video transformers with spatial-temporal token selection. In *European Conference on Computer Vision*, pages 69–86. Springer, 2022.
- [66] Yingquan Wang, Pingping Zhang, Shang Gao, Xia Geng, Hu Lu, and Dong Wang. Pyramid spatial-temporal aggregation for video-based person re-identification. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12026–12035, 2021.
- [67] Siyuan Wei, Tianzhu Ye, Shen Zhang, Yao Tang, and Jiajun Liang. Joint token pruning and squeezing towards more aggressive compression of vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2092–2101, 2023.
- [68] Thomas Wiegand, Gary J Sullivan, Gisle Bjontegaard, and Ajay Luthra. Overview of the h. 264/avc video coding standard. *IEEE Transactions on circuits and systems for video technology*, 13(7):560–576, 2003.
- [69] Jinlin Wu, Lingxiao He, Wu Liu, Yang Yang, Zhen Lei, Tao Mei, and Stan Z Li. Cavit: Contextual alignment vision transformer for video object re-identification. In *European Conference on Computer Vision*, pages 549–566. Springer, 2022.
- [70] Kan Wu, Jinnian Zhang, Houwen Peng, Mengchen Liu, Bin Xiao, Jianlong Fu, and Lu Yuan. Tinyvit: Fast pretraining distillation for small vision transformers. In *European Conference on Computer Vision*, pages 68–85. Springer, 2022.
- [71] Lin Yuanbo Wu, Lingqiao Liu, Yang Wang, Zheng Zhang, Farid Bousaid, Mohammed Bennamoun, and Xianghua Xie. Learning resolution-adaptive representations for cross-resolution person re-identification. *IEEE Transactions on Image Processing*, 2023.
- [72] Wei Wu, Jiawei Liu, Kecheng Zheng, Qibin Sun, and Zheng-Jun Zha. Temporal complementarity-guided reinforcement learning for image-to-video person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7319–7328, 2022.
- [73] Yiming Wu, Omar El Farouk Bourahla, Xi Li, Fei Wu, Qi Tian, and Xue Zhou. Adaptive graph representation learning for video person re-identification. *IEEE Transactions on Image Processing*, 29:8821–8830, 2020.
- [74] Enze Xie, Wenhao Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34:12077–12090, 2021.
- [75] Yichao Yan, Jie Qin, Jiabin Chen, Li Liu, Fan Zhu, Ying Tai, and Ling Shao. Learning multi-granular hypergraphs for video-based person re-identification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2899–2908, 2020.
- [76] Jinrui Yang, Wei-Shi Zheng, Qize Yang, Ying-Cong Chen, and Qi Tian. Spatial-temporal graph convolutional network for video-based person re-identification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3289–3299, 2020.
- [77] Xi Yang, Xian Wang, Liangchen Liu, Nannan Wang, and Xinbo Gao. Sfte: A comprehensive video-based person re-identification network based on spatio-temporal feature enhancement. *IEEE Transactions on Multimedia*, 2024.
- [78] Haoran You, Huihong Shi, Yipin Guo, and Yingyan Lin. Shiftaddvit: Mixture of multiplication primitives towards efficient vision transformer. *Advances in Neural Information Processing Systems*, 36, 2024.
- [79] Haoran You, Yunyang Xiong, Xiaoliang Dai, Bichen Wu, Peizhao

- Zhang, Haoqi Fan, Peter Vajda, and Yingyan Celine Lin. Castling-vit: Compressing self-attention via switching towards linear-angular attention at vision transformer inference. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14431–14442, 2023.
- [80] Chenyang Yu, Xuehu Liu, Yingquan Wang, Pingping Zhang, and Huchuan Lu. Tf-clip: Learning text-free clip for video-based person re-identification. *Proceedings of the AAAI conference on artificial intelligence*, 2024.
- [81] Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan. Metaformer is actually what you need for vision. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10819–10829, 2022.
- [82] Zhihang Yuan, Chenhao Xue, Yiqi Chen, Qiang Wu, and Guangyu Sun. P4vit: Post-training quantization for vision transformers with twin uniform quantization. In *European conference on computer vision*, pages 191–207. Springer, 2022.
- [83] Ying Zang, Runlong Cao, Chenglong Fu, Didi Zhu, Min Zhang, Wenjun Hu, Lanyun Zhu, and Tianrun Chen. Resmatch: Referring expression segmentation in a semi-supervised manner. *Information Sciences*, 694:121709, 2025.
- [84] Jinnian Zhang, Houwen Peng, Kan Wu, Mengchen Liu, Bin Xiao, Jianlong Fu, and Lu Yuan. Minivit: Compressing vision transformers with weight multiplexing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12145–12154, 2022.
- [85] Pingping Zhang, Yuhao Wang, Yang Liu, Zhengzheng Tu, and Huchuan Lu. Magic tokens: Select diverse tokens for multi-modal object re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17117–17126, 2024.
- [86] Wei Zhang, Xiaodong Yu, and Xuanyu He. Learning bidirectional temporal cues for video-based person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(10):2768–2776, 2017.
- [87] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, and Zhibo Chen. Multi-granularity reference-aided attentive feature aggregation for video-based person re-identification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10407–10416, 2020.
- [88] Linqing Zhao, Yi Wei, Jiabin Li, Jie Zhou, and Jiwen Lu. Structure-aware cross-modal transformer for depth completion. *IEEE Transactions on Image Processing*, 2024.
- [89] Yiru Zhao, Xu Shen, Zhongming Jin, Hongtao Lu, and Xian-sheng Hua. Attribute-driven feature disentangling and temporal aggregation for video person re-identification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4913–4922, 2019.
- [90] Ce Zheng, Sijie Zhu, Matias Mendieta, Taojiannan Yang, Chen Chen, and Zhengming Ding. 3d human pose estimation with spatial and temporal transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11656–11665, 2021.
- [91] Liang Zheng, Yi Yang, and Alexander G Hauptmann. Person re-identification: Past, present and future. *arXiv preprint arXiv:1610.02984*, 2016.
- [92] Zhedong Zheng, Xiaodong Yang, Zhiding Yu, Liang Zheng, Yi Yang, and Jan Kautz. Joint discriminative and generative learning for person re-identification. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2138–2147, 2019.
- [93] Kaiyang Zhou, Yongxin Yang, Andrea Cavallaro, and Tao Xiang. Omni-scale feature learning for person re-identification. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3702–3712, 2019.
- [94] Zhen Zhou, Yan Huang, Wei Wang, Liang Wang, and Tieniu Tan. See the forest for the trees: Joint spatial and temporal recurrent neural networks for video-based person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4747–4756, 2017.
- [95] Lanyun Zhu, Tianrun Chen, Deyi Ji, Jieping Ye, and Jun Liu. Llaf: When large language models meet few-shot segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3065–3075, 2024.
- [96] Lanyun Zhu, Tianrun Chen, Jianxiang Yin, Simon See, and Jun Liu. Continual semantic segmentation with automatic memory sample selection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3082–3092, 2023.
- [97] Lanyun Zhu, Tianrun Chen, Jianxiang Yin, Simon See, and Jun Liu. Learning gabor texture features for fine-grained recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1621–1631, 2023.
- [98] Lanyun Zhu, Tianrun Chen, Jianxiang Yin, Simon See, and Jun Liu. Addressing background context bias in few-shot segmentation through iterative modulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3370–3379, 2024.
- [99] Lanyun Zhu, Deyi Ji, Tianrun Chen, Peng Xu, Jieping Ye, and Jun Liu. Ibd: Alleviating hallucinations in large vision-language models via image-biased decoding. *arXiv preprint arXiv:2402.18476*, 2024.
- [100] Lanyun Zhu, Deyi Ji, Shiping Zhu, Weihao Gan, Wei Wu, and Junjie Yan. Learning statistical texture for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12537–12546, 2021.
- [101] Shiping Zhu, Chang Liu, and Ziyao Xu. High-definition video compression system based on perception guidance of salient information of a convolutional neural network and hevcc compression domain. *IEEE Transactions on Circuits and Systems for Video Technology*, 30(7):1946–1959, 2019.
- [102] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020.