# MSINet: Twins Contrastive Search of Multi-Scale Interaction for Object ReID

Jianyang Gu[1]    Kai Wang[2]    Hao Luo[3]    Chen Chen[4]    Wei Jiang[1*]
Yuqiang Fang[5]    Shanghang Zhang[6]    Yang You[2]    Jian Zhao[7*]

[1]Zhejiang University    [2]National University of Singapore    [3]Alibaba Group
[4]OPPO Research Institute    [5]Space Engineering University    [6]Peking University
[7]Institute of North Electronic Equipment

{gu_jianyang, jiangwei_zju}@zju.edu.cn zhaojian90@u.nus.edu

## Abstract

*Neural Architecture Search (NAS) has been increasingly appealing to the society of object Re-Identification (ReID), for that task-specific architectures significantly improve the retrieval performance. Previous works explore new optimizing targets and search spaces for NAS ReID, yet they neglect the difference of training schemes between image classification and ReID. In this work, we propose a novel Twins Contrastive Mechanism (TCM) to provide more appropriate supervision for ReID architecture search. TCM reduces the category overlaps between the training and validation data, and assists NAS in simulating real-world ReID training schemes. We then design a Multi-Scale Interaction (MSI) search space to search for rational interaction operations between multi-scale features. In addition, we introduce a Spatial Alignment Module (SAM) to further enhance the attention consistency confronted with images from different sources. Under the proposed NAS scheme, a specific architecture is automatically searched, named as MSINet. Extensive experiments demonstrate that our method surpasses state-of-the-art ReID methods on both in-domain and cross-domain scenarios. Source code available in https://github.com/vimar-gu/MSINet.*

## 1. Introduction

Object re-identification (Re-ID) aims at retrieving specific object instances across different views [39, 40, 57, 65, 70], which attracts much attention in computer vision community due to its wide-range applications. Previous works have achieved great progresses on both supervised [42, 49, 58] and unsupervised ReID tasks [17, 50, 78], most of which adopts backbone models originally designed for general image classification tasks [20, 52].

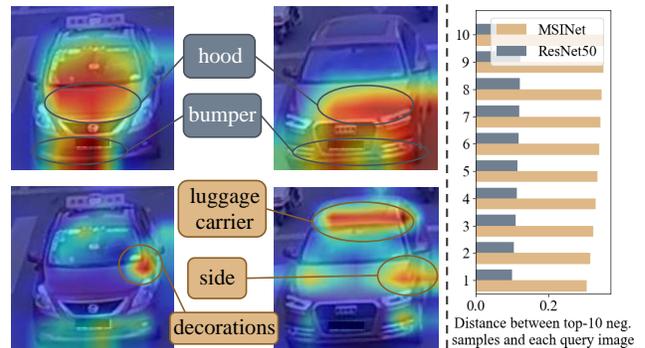Recent literature [64, 75] has shown that applying differ-



Figure 1. The left panel shows the example activation maps of ResNet50 (1st row) and MSINet (2nd row). The right panel shows the average distances between the most similar 10 negative samples and each query image at the inference. Best viewed in color.

ent architectures on ReID leads to large performance variations. Some works employ Neural Architecture Search (NAS) for ReID [28, 45]. The proposed optimizing targets and search spaces stably improve the model performance, yet the main search scheme still follows traditional NAS methods designed for general classification tasks [12, 36]. As an open-set task, ReID contains different categories in the training and validation sets [64, 71], while the two sets share exactly the same categories in standard classification tasks [10], which is also followed by traditional NAS methods. The incompatibility between search schemes and real-world training schemes makes the searched architecture sub-optimal for ReID. Moreover, ReID is required to distinguish more subtle distinctions among fine-grained instances compared with image-level classification [48, 63]. Some previous works [4, 44, 68, 75] have manifested that local perspectives and multi-scale features are discriminative for ReID. However, current utilizations of these features are mostly empirically designed, which can be more flexible according to the characteristics of different network layers.

In this work, we propose a novel NAS scheme aiming at addressing the aforementioned challenges. In order to simulate the real-world ReID training schemes, a Twins Contrastive Mechanism (TCM) is proposed to unbind the categories of the training and validation sets. An adjustable overlap ratio of categories builds up the compatibility between NAS and ReID, which provides more appropriate supervision for ReID architecture search. Moreover, to search for more rational utilizations of multi-scale features, we design a Multi-Scale Interaction (MSI) search space. The MSI space focuses on interaction operations between multi-scale features along the shallow and deep layers of the network, which guides the features to promote each other. Additionally, to further improve the generalization capability, we propose a Spatial Alignment Module (SAM) to enhance the attention consistency of the model confronted with images from different sources. With the above NAS scheme, we obtain a light-weight yet effective model architecture, denoted as Multi-Scale Interaction Net (MSINet).

We visualize the example activation maps of our proposed MSINet and ResNet50 [20] trained on VeRi-776 [38, 39] in Fig. 1. Compared to ResNet50, MSINet focuses on more unique distinctions with specific semantic information to recognize instances. Besides, MSINet largely increases the distance margin between query image and corresponding negative samples, reflecting extraordinary discriminative capability. Extensive experiments demonstrate that MSINet surpasses state-of-the-art (SOTA) ReID methods on both in-domain and cross-domain scenarios. Our source codes are available in the supplementary material.

Our contributions are summarized as follows:

- To the best of our knowledge, we are the first to build the NAS search scheme according to the real-world ReID training schemes, which provides more appropriate supervision for the ReID architecture search.

- We propose a novel search space based on the Multi-Scale Interaction (MSI) operations and a Spatial Alignment Module (SAM) to improve the model performance on in-domain and cross-domain scenarios.

- We construct a light-weight yet effective architecture for ReID tasks, denoted as MSINet. With only 2.3M parameters, MSINet surpasses ResNet50 [20] by 9% mAP on MSMT17 [60] and 16% mAP on MSMT17→Market-1501 [69].

## 2. Related Works

**Neural Architecture Search.** NAS has been increasingly appealing to the computer vision society, due to its automatic architecture designing characteristics. NAS methods can be roughly separated into four categories: reinforcement learning [1, 77], evolutionary algorithms [35, 46], gra-

dient desent [36, 43] and performance prediction [11, 34]. Liu *et al.* establish a differentiable architecture search (DARTS) method [36], which improves the practicability of NAS by a large extent. Some later works further improve the structure through sampling strategy [62], network pruning [2, 8], progressive learning [5], collaborative competition [7], *etc*. Most of NAS works focus on general image classification tasks, where the training and validation sets share the exact same categories. Following the setting, however, leads to incompatibility with the real-world training schemes of object ReID. In this work, we unbind the category bond between the two sets and propose a novel search scheme suitable for ReID.

**ReID Network Design.** Current ReID works mostly adopt backbones designed for image classification [20, 23, 47, 52]. Some works [14, 31, 67] design attention modules based on the common backbones to unearth their potential on distinguishing local distinctions. However, these methods usually lead to large calculation consumption.

There are also several works focusing on designing ReID-specific architectures. Li *et al.* present a Filter Pairing Neural Network to dynamically match patches in the feature maps [30]. Wang *et al.* separate and regroup the features of two samples with a WConv layer [56]. Guo *et al.* extract multi-scale features to directly evaluate the similarity between samples [18]. However, the siamese structure is inconvenient when conducting retrieval on large galleries. Zhou *et al.* aggregate multi-scale information to achieve high accuracy with small computing consumption [75]. Quan *et al.* introduce a part-aware module into the DARTS search space [36, 45]. Li *et al.* propose a new search space in regard to receptive field scales [28]. These methods have excellent performance on limited parameter scales, but fail to surpass those networks with complex structures. Different from previous works, we design a light-weight searching structure focusing on rational interaction operations between multi-scale features. The searched MSINet surpasses SOTA methods on both in-domain and cross-domain tasks.

## 3. Methods

Our goal is to construct an effective NAS scheme to search for a light-weight backbone architecture suitable for ReID tasks. Based on the training schemes of ReID, we propose a novel Twins Contrastive Mechanism to provide more appropriate supervision for the search process. Aiming at rational interaction between multi-scale features, we design a Multi-Scale Interaction search space. We further introduce a Spatial Alignment Module to improve the generalization capability with limited parameter growth.

### 3.1. Twins Contrastive Mechanism

NAS aims at automatically searching for the optimal network architecture for certain data. Inspired by [36], a basic
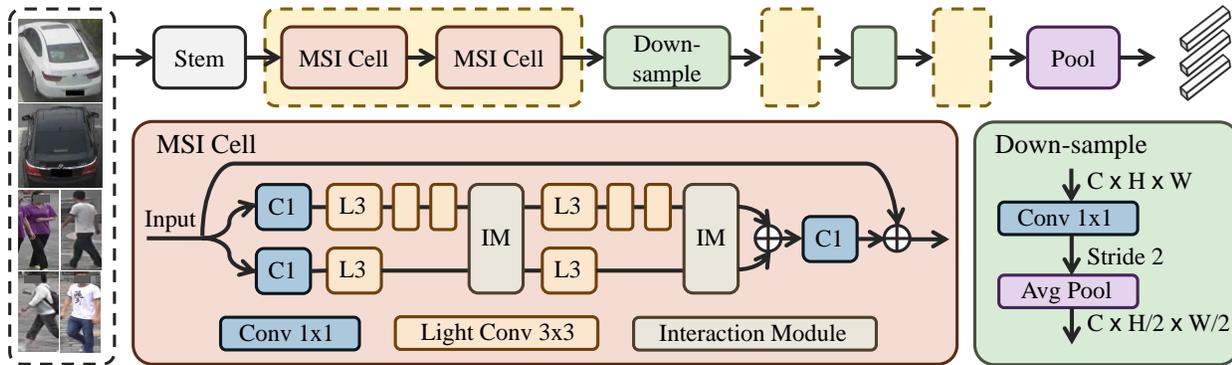
Figure 2. The model structure of the proposed MSINet. The input can be either *person* or *vehicle*. Inside a cell, the input is separated to two branches, with different receptive field scales. The interaction module exchanges information between two branches. Architecture search automatically select the most appropriate interaction for each cell.

differentiable architecture search scheme is established. We define the ordinary model parameters as $\omega$, and architecture parameters as $\alpha$. For network layer $i$ with a search space of $\mathcal{O}$, $\alpha_i$ controls the weight of each operation $o$ in the space. The features are parallelly passed through all the operations, and the final output is formulated by the softmax-weighted sum of operation outputs:

$$f(\mathbf{x}_i) = \sum_{o \in \mathcal{O}} \frac{\exp\{\alpha_i^o\}}{\sum_{o' \in \mathcal{O}} \exp\{\alpha_i^{o'}\}} \cdot o(\mathbf{x}_i). \quad (1)$$

The search process is conducted in an alternative manner. Training data is utilized to update the model parameters, and validation data is then employed to update the architecture parameters. For most NAS methods designed for image classification tasks, the training and validation data share exactly the same categories and a linear classification layer for loss calculation.

Different from standard image classification, as an open-set retrieval task, ReID has different categories in the training and validation sets. The incompatibility between search schemes and real-world training schemes might lead to sub-optimal searching results. Accordingly, we propose a novel Twins Contrastive Mechanism (TCM) for NAS ReID training. Specifically, we employ two independent auxiliary memories $\mathcal{C}_{tr}$ and $\mathcal{C}_{val}$ to store the embedded features of the training and validation data, respectively. The memories are initialized with the centroid features, which are calculated by averaging the features of each category. At each iteration, the training loss is first calculated with $\mathcal{C}_{tr}$ for model parameter updating. Given an embedded feature $\mathbf{f}$ with category label $j$, the contrastive classification loss is calculated with:

$$\mathcal{L}_{tr}^{cls} = -\log \frac{\exp(\mathbf{f} \cdot \mathbf{c}_{tr}^j/\tau)}{\sum_{n=0}^{N_{tr}^c} \exp(\mathbf{f} \cdot \mathbf{c}_{tr}^n/\tau)}, \quad (2)$$

where $\mathbf{c}_{tr}^n$ represents the memorized feature of category $n$, $N_{tr}^c$ stands for the total number of categories in the training set, and $\tau$ is the temperature parameter, which is set as 0.05 empirically [17]. After updating the model parameters, the embedded feature $\mathbf{f}$ with category label $j$ is integrated into the corresponding memorized feature $\mathbf{c}_{tr}^j$ by:

$$\mathbf{c}_{tr}^j \leftarrow \beta \mathbf{c}_{tr}^j + (1-\beta)\mathbf{f}, \quad (3)$$

where $\beta$ is set as 0.2 empirically [17]. Then the updated model is evaluated on the validation data to generate to validation loss with $\mathcal{C}_{val}$ replacing $\mathcal{C}_{tr}$ in Eq. 2. The architecture parameter is then updated with the validation loss to finish an iteration.

As the loss calculation does not rely on the linear classification layer, the categories of the training and validation sets are unbound. We are able to dynamically adjust the category overlap ratio in these two sets. The advantages of a proper overlap ratio are summarized as two folds. Firstly, TCM better simulates the real-world training of ReID and helps the model focus on truly discriminative distinctions. The differences between the training and validation data improves the generalization capability of the model. Secondly, a relatively small proportion of overlapped categories stabilizes the architecture parameter update through a consistent optimizing target with the model parameter update.

## 3.2. Multi-Scale Interaction Space

Although the local perspective and multi-scale features have already been investigated in previous ReID works [4, 28, 44, 51, 68, 72, 75, 76], the utilization of these information is mainly empirically designed aggregation, which is monotonous and restrained. We argue that on the one hand, the rational utilization of multi-scale features should be dynamically adjusted along the shallow and deep layers of the network. On the other hand, introducing interaction other than aggregation creates direct information exchange, and
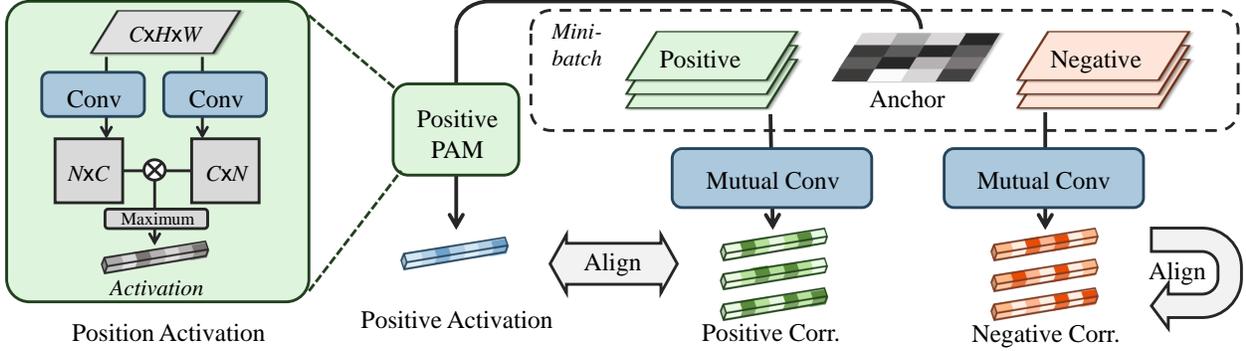
Figure 3. The aligning pipeline of the proposed Spatial Alignment Module. The correlation activation vectors are calculated between the anchor feature and all features in a mini-batch. The positive vectors are aligned with the learnable self-activation, and the negative vectors are aligned with each other. The structure of PAM is shown on the left.

Table 1. The detailed interaction operation in the proposed MSINet architecture. N: None; E: Exchange; G: Channel Gate; A: Cross Attention.

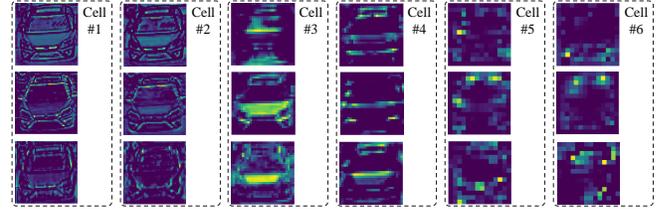| Cell #1 | | Cell #2 | | Cell #3 | | Cell #4 | | Cell #5 | | Cell #6 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| G | G | E | G | A | G | G | N | G | A | E | A |



Figure 4. Output feature maps of MSI cells at different layers.

makes fuller use of multi-scale features. Therefore, we propose a novel Multi-Scale Interaction (MSI) search space to establish a light-weight architecture suitable for ReID.

As shown in Fig. 2, the network is mainly grouped with MSI cells and down-sample blocks, which is generally consistent with OSNet [75]. In each cell, the input features are passed through two branches with different receptive field scales. To reduce the calculation burden of the network, for the layers inside each branch, we adopt the stack of 1×1 convolution and multiple depth-wise 3×3 convolution to implement specific scales. A scale ratio $\rho$ of 3:1 is selected for the two branches. These two branches do not share model parameters, except for the Interaction Modules (IM). IM introduces information exchange for the two branches. There are 4 operation options for the IM. With the two-branch input features defined as $(\mathbf{x}_1, \mathbf{x}_2)$, the operations can be formulated as:

**None.** None operation involves no parameters, and outputs exactly the input features $(\mathbf{x}_1, \mathbf{x}_2)$.

**Exchange.** Exchange acts as the strongest interaction among all options. It directly exchanges the features for the two branches and outputs $(\mathbf{x}_2, \mathbf{x}_1)$. Exchange contains no extra parameters, as well.

**Channel Gate.** Channel gate introduces a Multi-Layer Perceptron (MLP) to generate a channel-wise attention gate [61, 75] as:

$$G(\mathbf{x}) = \sigma(MLP(\mathbf{x})), \qquad (4)$$

and returns $(G(\mathbf{x}_1) \cdot \mathbf{x}_1, G(\mathbf{x}_2) \cdot \mathbf{x}_2)$. The MLP is composed of 2 fully connected layers and its parameters are shared for both branches. Thereby it achieves interaction by jointly screening discriminative feature channels.

**Cross Attention.** Traditional channel attention module calculates the channel correlation inside a single feature map [15]. The original feature map $\mathbf{x} \in \mathcal{R}^{C \times H \times W}$ is firstly reshaped into the *query* feature $\tilde{\mathbf{x}} \in \mathcal{R}^{C \times N}$, where $N = H \times W$. Then the correlation activation is calculated by performing a matrix multiplication between the *query* feature $\tilde{\mathbf{x}}$ and the *key* feature $\tilde{\mathbf{x}}^\top$. We propose to exchange the *keys* of the two branches to explicitly calculate the correlation between each other. The correlation activation is then transformed to a mask, and is added up to the original features with a learnable proportion.

After interaction, the multi-scale branches are fused through a sum operation. It is worth noting that the extra parameters brought by multiple interaction options are limited, which enables searching for each cell along the whole network independently. At the beginning of the network, we employ the same stem module as that in OSNet [75], containing a $7 \times 7$ convolutional layer and a $3 \times 3$ max pooling with a stride of 2. After the searching process, the interaction operation $o$ with the largest weight $\alpha_i^o$ at each layer is reserved to form the searched architecture.

After searching the architecture, the model is validated on various Re-ID tasks. The training is constrained by the

4

Table 2. Supervised performance on object ReID datsets. The results in the top part are trained from scratch, and those in the bottom part are pre-trained on ImageNet in advance. As the compared methods are originally proposed for person ReID, we reproduce the results in vehicle datasets. $^*$ indicates that the results of person ReID are reproduced by us. The evaluation results of architecture searched on VR can be found in the supplementary material.

| Method | Params | Inference Time | M | | MS | | VR | | VID | | MS→M | | VR→VID | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | R-1↑ | mAP↑ | R-1↑ | mAP↑ | R-1↑ | mAP↑ | R-1↑ | R-5↑ | R-1↑ | mAP↑ | R-1↑ | R-5↑ |
| ResNet50$^*$ [42] | ∼24M | 1× | 85.7 | 68.3 | 48.0 | 25.7 | 92.8 | 69.9 | 70.6 | 76.6 | - | - | | |
| OSNet [75] | 2.2M | 0.79× | 93.6 | 81.0 | 71.0 | 43.3 | 95.4 | 72.8 | 76.0 | 88.7 | - | - | - | - |
| CDNet [28] | 1.8M | 0.67× | 93.7 | 83.7 | 73.7 | 48.5 | 94.3 | 73.0 | 74.5 | 88.8 | - | - | - | - |
| MSINet | 2.3M | 0.71× | **94.6** | **87.0** | **76.0** | **52.5** | **95.9** | **75.0** | **76.5** | **89.8** | - | - | - | - |
| ResNet50$^*$ [42] | ∼24M | 1× | 94.5 | 85.9 | 75.5 | 50.4 | 94.5 | 73.6 | 76.5 | 89.9 | 58.8 | 31.8 | 42.8 | 61.9 |
| OSNet [75] | 2.2M | 0.79× | 94.8 | 84.9 | 78.7 | 52.9 | 95.5 | 76.4 | 76.0 | 88.6 | 66.6 | 37.5 | 46.5 | 63.1 |
| CDNet [28] | 1.8M | 0.67× | 95.1 | 86.0 | 78.9 | 54.7 | - | - | - | - | - | - | - | - |
| MSINet | 2.3M | 0.71× | 95.3 | 89.6 | **81.0** | **59.6** | **96.8** | 78.8 | 77.9 | 91.7 | 74.9 | 46.2 | 48.0 | 65.6 |
| MSINet-SAM | 2.4M | 0.71× | **95.5** | **89.9** | 80.7 | 59.5 | 96.7 | **79.0** | **78.0** | **91.9** | 76.3 | 48.4 | 49.0 | 66.8 |

classification id loss and the triplet loss, formulated by:

$$\mathcal{L}_{id} = \frac{1}{N} \sum_{i=1}^{N} - \log \left( \frac{\exp \mathbf{W}_i^\top \mathbf{f}_i}{\sum_j \exp \mathbf{W}_j^\top \mathbf{f}_i} \right), \qquad (5)$$

where $\mathbf{f}_i$ is a feature vector, the corresponding classifier weight of which is $\mathbf{W}_i$, and

$$\mathcal{L}_{tri} = \left[ \mathcal{D}(\mathbf{f}_a, \mathbf{f}_p) - \mathcal{D}(\mathbf{f}_a, \mathbf{f}_n) + m \right]_+, \qquad (6)$$

where $\mathbf{f}_a$, $\mathbf{f}_p$, $\mathbf{f}_n$ are the embedded features for the anchor, the hardest positive and negative samples in a mini-batch, $\mathcal{D}(\cdot, \cdot)$ is the Euclidean distance, $m$ is the margin parameter, and $[\cdot]_+$ is the $\max(\cdot, 0)$ function.

### 3.3. Spatial Alignment Module

The retrieval precision of object ReID tasks are largely affected by the variation of appearances such as poses, illumination and occlusion when the camera conditions change. In order that the model correctly and consistently focuses on the discriminative spatial positions, we design a Spatial Alignment Module (SAM) to explicitly align the spatial attention between images, as shown in Fig. 3.

Specifically, we first calculate the position-wise correlation activation map $\mathbf{A}$ between the feature maps in a mini-batch. The activation between sample $i$ and $j$ can be formulated as: $\mathbf{A}(i,j) = \tilde{\mathbf{x}}_j^\top \times \tilde{\mathbf{x}}_i$, where $\tilde{\mathbf{x}} \in \mathcal{R}^{C \times N}$ is reshaped from the original feature $\mathbf{x} \in \mathcal{R}^{C \times H \times W}$. Then we take the maximum activation for each position of sample $i$ as:

$$\mathbf{a}(i,j) = \max_{dim=1} \mathbf{A}(i,j). \qquad (7)$$

The above process is denoted as "Mutual Conv" in Fig. 3. We evaluate the consistency between activation vectors with cosine similarity. For negative samples specifically, there

can be many different hints for recognition, some of which might be inappropriate, such as the backgrounds. By aligning all the correlations for sample $i$, we hope that the network can correct some attention bias and consistently focus on discriminative positions.

However, through aligning positive sample pairs, the ID-related features are expected to be emphasized, which cannot be achieved by aligning negative pairs. Therefore, we introduce an extra position activation module (PAM) to generate supervision for the alignment between positive pairs. The spatial alignment loss is formulated as:

$$\mathcal{L}_{sa}(i) = \frac{1}{N_+} \sum_{p \in \mathcal{I}_+} (1 - S(\hat{\mathbf{a}}(i), \mathbf{a}(i,p))) + \frac{1}{N_-} \sum_{n_1, n_2 \in \mathcal{I}_-} (1 - S(\mathbf{a}(i,n_1), \mathbf{a}(i,n_2))), \qquad (8)$$

where $\mathcal{I}_+$ contains positive indices for sample $i$, the total number of which is $N_+$, and vice versa. $\hat{\mathbf{a}}(i)$ stands for the generated activation vector for positive sample alignment, and $S(\cdot, \cdot)$ is the cosine similarity.

## 4. Experiments

### 4.1. Datasets and Evaluation Metrics

Our proposed method is evaluated on two person ReID datasets Market-1501 [69], MSMT17 [60], and two vehicle ReID datasets VeRi-776 [38, 39] and VehicleID [37]. For simplicity, the four datasets are denoted as M, MS, VR and VID in the following sections, respectively. Evaluation metrics include Cumulative Matching Characteristic (CMC) and mean average precision (mAP), which are commonly utilized on ReID tasks.

Table 3. Supervised performance comparison between MSINet and SOTA methods on M and MS datasets.

| Method | M | | MS | |
| --- | --- | --- | --- | --- |
| | R-1↑ | mAP↑ | R-1↑ | mAP↑ |
| PCB [51] | 93.8 | 81.6 | 68.2 | 40.4 |
| MGN [55] | 95.7 | 86.9 | 76.9 | 52.1 |
| OSNet [75] | 93.6 | 81.0 | 71.0 | 43.3 |
| IANet [22] | 94.4 | 83.1 | 75.5 | 46.8 |
| DGNet [73] | 94.8 | 86.0 | 77.2 | 52.3 |
| Auto-ReID [45] | 94.5 | 85.1 | - | - |
| SAN [26] | **96.1** | 88.0 | 79.2 | 55.7 |
| CDNet [28] | 95.1 | 86.0 | 78.9 | 54.7 |
| BAT-Net [14] | 95.1 | 87.4 | 79.5 | 56.8 |
| SFT [41] | 94.1 | 87.5 | 79.0 | 58.3 |
| CTF [66] | 94.8 | 87.7 | - | - |
| RGA-SC [67] | **96.1** | 88.4 | 80.3 | 57.5 |
| **MSINet** | 95.3 | **89.6** | **81.0** | **59.6** |

Table 4. Unsupervised performance applying MSINet to SOTA methods for USL on M and UDA on M→MS.

| Method | M | | M→MS | |
| --- | --- | --- | --- | --- |
| | R-1↑ | mAP↑ | R-1↑ | mAP↑ |
| MMCL [54] | 80.3 | 45.5 | 40.8 | 15.1 |
| MMT [16] | - | - | 50.1 | 24.0 |
| JVTC+ [29] | 79.5 | 47.5 | 48.6 | 25.1 |
| CycAs [59] | 84.8 | 64.8 | - | - |
| GCL [3] | 87.3 | 66.8 | 51.1 | 27.0 |
| MPRD [24] | 83.0 | 51.0 | - | - |
| SpCL [17] | 88.1 | 73.1 | 53.7 | 26.8 |
| GS [19] | 92.3 | 79.2 | - | - |
| **GS+MSINet** | 91.7 | 81.5 | - | - |
| HDCRL [6] | 92.4 | 81.7 | - | - |
| **HDCRL+MSINet** | **92.9** | **82.7** | - | - |
| IDM [9] | - | - | 61.3 | 33.5 |
| **IDM+MSINet** | - | - | **66.0** | **37.8** |

## 4.2. Architecture Search

We conduct the searching process on MSMT17. SGD is adopted for model parameter update with an initial learning rate of 0.025. The model is trained for 350 epochs in total. We adopt a warm-up strategy for the first 10 epochs. Then the learning rate is decayed by 0.1 at 150, 225 and 300 epochs, respectively. Adam [27] is adopted for the architecture parameter update with an initial learning rate of 0.002. The learning rate is decayed at the same pace. The images are reshaped to 256×128 for person and 256×256 for vehicles. Data augmentation includes random flip, random crop and random erasing [74]. The searched architecture is presented in Tab. 1. The "MSINet" in the following experiment sections refers to this architecture.

We visualize the feature maps extracted by each MSI cell in Fig. 4. At the shallow layers of the network, the kernels mainly focus on overall contour information. Channel gate helps to filter out inferior information, such as the background. As we approach deeper layers, the extracted features each have specific semantic information, where cross attention is more likely to be selected for the interaction. It indicates that cross attention is more rational for exchanging high-level semantic information.

## 4.3. Comparison with Other Backbones

We first compare our proposed MSINet with ResNet50 and recent proposed light-weight backbones in both in-domain and cross-domain ReID scenarios.

**In-Domain Tasks.** We adopt a two-group supervised evaluation scheme similar to that in [28, 75]: training from scratch and fine-tuning ImageNet [10] pre-trained models. The training parameters for both schemes are kept the same as that in architecture search, except for an initial learning rate of 0.065. Triplet loss and cross entropy loss are adopted

for the parameter update. The margin $m$ in Eq. 6 is set as 0.3. [28] adopts an FBLNeck. We also employ the same structure. The results are shown in Tab. 2.

ResNet50 is the most commonly utilized backbone network in ReID tasks, yet holds the worst performance. Moreover, ResNet50 largely depends on ImageNet pre-training, while MSINet without pre-training has already surpassed pre-trained ResNet50 on all metrics. Compared with the other datasets, MS contains more variations on illumination, background and camera pose, and brings a large performance gap between ResNet50 and other methods. It also validates the inadequacy of image classification networks on ReID tasks. OSNet [75] and CDNet [28] are recently proposed architectures designed specifically for ReID tasks. Both architectures focus on fusing multi-scale features to better suit ReID. CDNet employs a traditional NAS scheme to search for the proper receptive field scales for each cell. MSINet fixes the receptive field scale and instead selects optimal interaction operations inside each cell. With only a bit more parameters, MSINet surpasses all the other backbones by a large margin.

**Cross-Domain Tasks.** Cross-domain experiments verify the generalization capability of the model. Following previous domain generalizable ReID works [25, 33], data augmentation is adjusted to random flip, random crop and color jittering. The model is pre-trained and fine-tuned for 250 epochs to avoid over-fitting. The other settings are kept the same as supervision scenes. With no present pre-trained models for CDNet [28], it is excluded from this section.

Tab. 2 shows that ResNet50 can be easily interfered by different image styles confronted with new image domains. OSNet learns multi-scale features with specific semantic information for ReID, which is domain invariant to some extent. Our proposed search scheme also takes into ac-
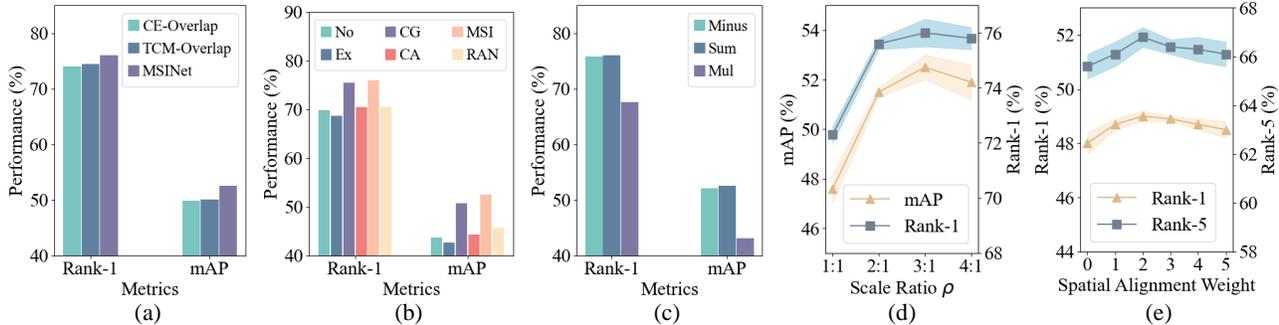
Figure 5. Ablation studies on (a) search scheme; (b) architecture; (c) multi-scale aggregation; (d) scale ratio $\rho$; (e) spatial alignment weight.

Table 5. The effectiveness of each components in SAM. Both in-domain performance on VR and cross-domain performance on VR→VID is evaluated.

| Model | Pos. | Neg. | Align | VR | | VR→VID | |
|---|---|---|---|---|---|---|---|
| | | | | R-1↑ | mAP↑ | R-1↑ | R-5↑ |
| MSINet | | | - | 96.7 | 78.5 | 48.0 | 65.6 |
| | ✓ | | Self | 96.7 | 78.1 | 48.4 | 66.0 |
| | | ✓ | Self | 96.7 | 78.4 | 48.5 | 66.2 |
| | ✓ | ✓ | Unified | 96.6 | 78.3 | 48.3 | 65.8 |
| | ✓ | ✓ | Separated | **96.8** | 78.6 | 48.7 | 66.4 |
| | ✓ | ✓ | PAM-Self | 96.7 | **79.0** | **49.0** | **66.8** |
| OSNet | | | None | 95.5 | 76.4 | 46.5 | 63.1 |
| | ✓ | ✓ | PAM-Self | 95.9 | 76.3 | 47.5 | 63.3 |

count the generalization capability of the model. By partly separating the categories for training and validation sets, the searched interaction operations generalize well confronted with new image domains. Except for discrimination, MSINet also surpasses the other backbones on cross-domain tasks by a large margin with faster inference speed.

Additionally, we introduce SAM into the model, which aligns the spatial correlations between images. A weighted sum of ReID loss and spatial alignment loss is utilized when training the network with SAM. The weight of spatial alignment loss is set as $\lambda_{sa} = 2.0$. Without extra inference consumption or damages on the supervised performance, SAM further boosts the generalization capability of MSINet.

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

Tab. 3 further illustrates the supervision performance comparison of our proposed MSINet with the SOTA methods on M and MS datasets. With much less parameters than most of the compared methods, MSINet achieves a retrieval accuracy comparable to that of more complicated ones. Auto-ReID [45] first designs a NAS scheme for ReID, yet the DARTS-style architecture contains 13M parameters. RGA-SC [67] carefully designs a relation-aware global at-

tention module. MSINet achieves even higher performance with less training consumption, which validates the superiority of selecting rational interaction.

We also evaluate the model performance replacing the backbone network from ResNet50 to MSINet for SOTA unsupervised ReID methods in Tab. 4. For purely unsupervised learning (USL) method GS [19] on M dataset, MSINet performs slightly lower on rank-1, yet has a large superiority on mAP. For HDCRL [6], MSINet shows obvious superiority over ResNet50. For unsupervised domain adaptation (UDA) method IDM [9] on M→MS task, MSINet surpasses ResNet50 by a large margin, which further proves that the TCM brings outstanding generalization capability to the searched architecture.

### 4.5. Ablation Studies

**Effectiveness of Architecture Search.** To verify the effectiveness, we conduct supervised training on MS with different search schemes in Fig. 5 (a). Under the standard classification scheme ("CE Overlap"), the searched model performs poor. Replacing the cross entropy loss to contrastive loss ("TCM Overlap") only brings a slight improvement. The complete TCM framework unbounds the categories between the training and validation set, and thereby improves the performance by a large margin.

In Fig. 5 (b) we compare the performance of different architectures. Firstly, we validate 4 models each with a unique interaction operation from the 4 options in the search space. None and Exchange, with no trainable parameters, achieve poor performance. Channel gate introduces channel-wise attention, whose model performs the best among 4 options. Cross attention exchanges the key features for the two branches. Over-frequent exchange interferes the ordinary feature extraction and degrades the network performance. Through appropriately arranging interaction operations along the architecture, MSINet surpasses all the above 4 models. Random architecture, on the other hand, shows no rational appliance of interaction operations, which validates that the proposed search scheme helps find

Table 6. Supervised performance comparison between MSINet and Transformers on VR and MS datasets.

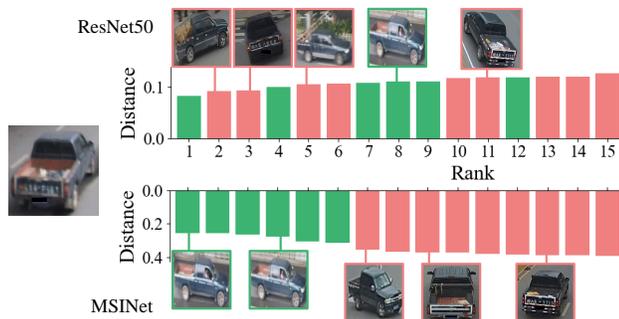| Method | Params | Inference Time | MS R-1↑ | MS mAP↑ | VR R-1↑ | VR mAP↑ |
|--------|--------|-------|------|-------|------|-------|
| DeiT-S [21] | ∼22M | 0.97x | 76.3 | 55.2 | 95.5 | 76.3 |
| DeiT-B [21] | ∼86M | 1.79x | **81.9** | **61.4** | 95.9 | 78.4 |
| ViT-B [21] | ∼86M | 1.79x | 81.8 | 61.0 | 96.5 | 78.2 |
| **MSINet** | 2.3M | 0.71x | 81.0 | 59.6 | **96.8** | **78.8** |



Figure 6. Example top-15 retrieved sequences comparison on VR. Appearance differences caused by variant camera conditions are well addressed by the proposed MSINet. Visualization of person ReID can be found in the supplementary material.

suitable architectures for ReID.

**Effectiveness of Spatial Alignment Module.** We validate the effectiveness of each components of SAM on the VR→VID cross-domain experiment in Tab. 5. Firstly, we introduce the spatial alignment for positive and negative sample pairs, respectively. Each of them brings certain performance improvements. However, a unified alignment for all sample pairs damages ID-related features and degrades the performance instead. Therefore, we separate the alignment of positive and negative samples, which retains some discriminative features and integrates the effect of both aspects. The extra PAM for positive sample alignment further guarantees the focus on ID-related positions and achieves the best performance. We also conduct in-domain experiment on VR to prove that SAM improves the generalization capability without sacrificing the supervision performance. Adding SAM to OSNet receives similar results, which validates the universality of SAM.

**Fusing Operation.** After interaction, the multi-scale features are fused by sum operation. We investigate several fusing options on MS training from scratch in Fig. 5 (c). Subtracting ("Minus") a branch from the other leads to about the same results as "Sum" while multiplication ("Mul") performs poorly.

**Comparison with Transformer.** Transformer, as a new architecture, has recently been continuously making progresses in many computer vision domains [13, 53], includ-

ing ReID [21, 32]. We compare the model performance with some baseline Transformer models in Tab. 6. DeiT-B and ViT-B [21] achieves higher performance on MSMT17, with much larger calculation burden compared with our proposed MSINet. On VeRi-776, MSINet surpasses all the baseline Transformer methods. It proves that rational interaction operations between multi-scale features are capable to assist light-weighted pure-CNN models to obtain comparable performance with complex Transformers.

**Parameter Analysis.** Firstly, we study the influence of different receptive field scale ratios $\rho$ inside an MSI cell on MS training from scratch in Fig. 5 (d). Introducing scale differences between branches improves the model performance significantly, and subsequent increases brings more modest impacts. Considering both parameter scales and model performance, the ratio of 3:1 is selected for MSINet.

Secondly, the model performance fluctuation influenced by spatial alignment weight is visualized in Fig. 5 (e). The experiment is conducted on the VR→VID cross-domain scenario. Employing the alignment generally makes a positive impact on the generalization capability of the model. The optimal loss weight $\lambda_{sa}$ locates at 2.0.

**Visualization Results.** We visualize the top-15 retrieved sequences and the corresponding distances from an query image on VR in Fig. 6. By comparison, ResNet50 mainly focuses on general appearance features, where the top-rated negative samples share similar car bodies. MSINet, oppositely, concentrates on discriminative distinctions, an empty car hopper in this case, and creates an evident distance gap between positive and negative samples. More details can be seen in the supplementation material.

## 5. Conclusion

In this paper, we design a Twins Contrastive Mechanism for NAS to build the compatibility with ReID. The task-specific search scheme provides the searching process with more appropriate supervision. A Multi-Scale Interaction search space is proposed to establish rational and flexible utilization of multi-scale features. With a Spatial Alignment Module, our proposed MSINet achieves SOTA performance on both supervision and cross-domain scenarios with limited parameter amount. We hope the proposed approach could inspire more works focusing on designing network architectures suitable for ReID tasks.

## 6. Acknowledgement

# References

[1] Irwan Bello, Barret Zoph, Vijay Vasudevan, and Quoc V Le. Neural optimizer search with reinforcement learning. In *ICML*, pages 459–468, 2017. 2

[2] Han Cai, Ligeng Zhu, and Song Han. Proxylessnas: Direct neural architecture search on target task and hardware. In *ICLR*, 2018. 2

[3] Hao Chen, Yaohui Wang, Benoit Lagadec, Antitza Dantcheva, and Francois Bremond. Joint generative and contrastive learning for unsupervised person re-identification. In *CVPR*, pages 2004–2013, 2021. 6

[4] Tsai-Shien Chen, Chih-Ting Liu, Chih-Wei Wu, and Shao-Yi Chien. Orientation-aware vehicle re-identification with semantics-guided part attention network. In *ECCV*, pages 330–346, 2020. 1, 3

[5] Xin Chen, Lingxi Xie, Jun Wu, and Qi Tian. Progressive differentiable architecture search: Bridging the depth gap between search and evaluation. In *ICCV*, pages 1294–1303, 2019. 2

[6] De Cheng, Jingyu Zhou, Nannan Wang, and Xinbo Gao. Hybrid dynamic contrast and probability distillation for unsupervised person re-id. *IEEE Transactions on Image Processing*, 31:3334–3346, 2022. 6, 7

[7] Xiangxiang Chu, Tianbao Zhou, Bo Zhang, and Jixiang Li. Fair darts: Eliminating unfair advantages in differentiable architecture search. In *ECCV*, pages 465–480, 2020. 2

[8] Xiyang Dai, Dongdong Chen, Mengchen Liu, Yinpeng Chen, and Lu Yuan. Da-nas: Data adapted pruning for efficient neural architecture search. In *ECCV*, pages 584–600, 2020. 2

[9] Yongxing Dai, Jun Liu, Yifan Sun, Zekun Tong, Chi Zhang, and Ling-Yu Duan. Idm: An intermediate domain module for domain adaptive person re-id. In *ICCV*, pages 11864–11874, 2021. 6, 7

[10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009. 1, 6

[11] Tobias Domhan, Jost Tobias Springenberg, and Frank Hutter. Speeding up automatic hyperparameter optimization of deep neural networks by extrapolation of learning curves. In *IJCAI*, 2015. 2

[12] Xuanyi Dong and Yi Yang. Searching for a robust neural architecture in four gpu hours. In *CVPR*, pages 1761–1770, 2019. 1

[13] 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. In *ICLR*, 2020. 8

[14] Pengfei Fang, Jieming Zhou, Soumava Kumar Roy, Lars Petersson, and Mehrtash Harandi. Bilinear attention networks for person retrieval. In *ICCV*, pages 8030–8039, 2019. 2, 6

[15] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. In *CVPR*, pages 3146–3154, 2019. 4

[16] Yixiao Ge, Dapeng Chen, and Hongsheng Li. Mutual mean-teaching: Pseudo label refinery for unsupervised domain adaptation on person re-identification. In *ICLR*, 2019. 6

[17] Yixiao Ge, Dapeng Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Self-paced contrastive learning with hybrid memory for domain adaptive object re-id. In *NeurIPS*, 2020. 1, 3, 6

[18] Yiluan Guo and Ngai-Man Cheung. Efficient and deep person re-identification using multi-level similarity. In *CVPR*, pages 2335–2344, 2018. 2

[19] Xumeng Han, Xuehui Yu, Nan Jiang, Guorong Li, Jian Zhao, Qixiang Ye, and Zhenjun Han. Group sampling for unsupervised person re-identification. *arXiv preprint arXiv:2107.03024*, 2021. 6, 7

[20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 1, 2, 12

[21] Shuting He, Hao Luo, Pichao Wang, Fan Wang, Hao Li, and Wei Jiang. Transreid: Transformer-based object re-identification. In *ICCV*, pages 15013–15022, 2021. 8

[22] Ruibing Hou, Bingpeng Ma, Hong Chang, Xinqian Gu, Shiguang Shan, and Xilin Chen. Interaction-and-aggregation network for person re-identification. In *CVPR*, pages 9317–9326, 2019. 6

[23] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *CVPR*, pages 4700–4708, 2017. 2

[24] Haoxuanye Ji, Le Wang, Sanping Zhou, Wei Tang, Nanning Zheng, and Gang Hua. Meta pairwise relationship distillation for unsupervised person re-identification. In *ICCV*, pages 3661–3670, 2021. 6

[25] Xin Jin, Cuiling Lan, Wenjun Zeng, Zhibo Chen, and Li Zhang. Style normalization and restitution for generalizable person re-identification. In *CVPR*, pages 3143–3152, 2020. 6

[26] Xin Jin, Cuiling Lan, Wenjun Zeng, Guoqiang Wei, and Zhibo Chen. Semantics-aligned representation learning for person re-identification. In *AAAI*, pages 11173–11180, 2020. 6

[27] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2014. 6

[28] Hanjun Li, Gaojie Wu, and Wei-Shi Zheng. Combined depth space based architecture search for person re-identification. In *CVPR*, pages 6729–6738, 2021. 1, 2, 3, 5, 6

[29] Jianing Li and Shiliang Zhang. Joint visual and temporal consistency for unsupervised domain adaptive person re-identification. In *ECCV*, pages 483–499, 2020. 6

[30] Wei Li, Rui Zhao, Tong Xiao, and Xiaogang Wang. Deepreid: Deep filter pairing neural network for person re-identification. In *CVPR*, pages 152–159, 2014. 2

[31] Wei Li, Xiatian Zhu, and Shaogang Gong. Harmonious attention network for person re-identification. In *CVPR*, pages 2285–2294, 2018. 2

[32] 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 *CVPR*, pages 2898–2907, 2021. 8

9

[33] Shengcai Liao and Ling Shao. Interpretable and generalizable person re-identification with query-adaptive convolution and temporal lifting. In *ECCV*, pages 456–474, 2020. 6

[34] Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search. In *ECCV*, pages 19–34, 2018. 2

[35] Hanxiao Liu, Karen Simonyan, Oriol Vinyals, Chrisantha Fernando, and Koray Kavukcuoglu. Hierarchical representations for efficient architecture search. In *ICLR*, 2018. 2

[36] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. In *ICLR*, 2018. 1, 2, 13

[37] Hongye Liu, Yonghong Tian, Yaowei Yang, Lu Pang, and Tiejun Huang. Deep relative distance learning: Tell the difference between similar vehicles. In *CVPR*, pages 2167–2175, 2016. 5

[38] Xinchen Liu, Wu Liu, Huadong Ma, and Huiyuan Fu. Large-scale vehicle re-identification in urban surveillance videos. In *ICME*, pages 1–6, 2016. 2, 5, 12

[39] Xinchen Liu, Wu Liu, Tao Mei, and Huadong Ma. A deep learning-based approach to progressive vehicle re-identification for urban surveillance. In *ECCV*, pages 869–884, 2016. 1, 2, 5, 12

[40] Xinchen Liu, Wu Liu, Tao Mei, and Huadong Ma. Provid: Progressive and multimodal vehicle reidentification for large-scale urban surveillance. *T-MM*, 20(3):645–658, 2017. 1

[41] Chuanchen Luo, Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. Spectral feature transformation for person re-identification. In *ICCV*, pages 4976–4985, 2019. 6

[42] Hao Luo, Wei Jiang, Youzhi Gu, Fuxu Liu, Xingyu Liao, Shenqi Lai, and Jianyang Gu. A strong baseline and batch normalization neck for deep person re-identification. *T-MM*, 22(10):2597–2609, 2019. 1, 5, 12, 13

[43] Renqian Luo, Fei Tian, Tao Qin, Enhong Chen, and Tie-Yan Liu. Neural architecture optimization. In *NeurIPS*, 2018. 2

[44] Xuelin Qian, Yanwei Fu, Yu-Gang Jiang, Tao Xiang, and Xiangyang Xue. Multi-scale deep learning architectures for person re-identification. In *ICCV*, pages 5399–5408, 2017. 1, 3

[45] Ruijie Quan, Xuanyi Dong, Yu Wu, Linchao Zhu, and Yi Yang. Auto-reid: Searching for a part-aware convnet for person re-identification. In *ICCV*, pages 3750–3759, 2019. 1, 2, 6, 7

[46] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V Le. Regularized evolution for image classifier architecture search. In *AAAI*, pages 4780–4789, 2019. 2

[47] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *CVPR*, pages 4510–4520, 2018. 2

[48] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, pages 815–823, 2015. 1

[49] Yantao Shen, Tong Xiao, Hongsheng Li, Shuai Yi, and Xiaogang Wang. Learning deep neural networks for vehicle re-id with visual-spatio-temporal path proposals. In *ICCV*, pages 1900–1909, 2017. 1

[50] Liangchen Song, Cheng Wang, Lefei Zhang, Bo Du, Qian Zhang, Chang Huang, and Xinggang Wang. Unsupervised domain adaptive re-identification: Theory and practice. *PR*, 102:107173, 2020. 1

[51] Yifan Sun, Liang Zheng, Yi Yang, Qi Tian, and Shengjin Wang. Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline). In *ECCV*, pages 480–496, 2018. 3, 6

[52] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, pages 2818–2826, 2016. 1, 2

[53] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, pages 10347–10357, 2021. 8

[54] Dongkai Wang and Shiliang Zhang. Unsupervised person re-identification via multi-label classification. In *CVPR*, pages 10981–10990, 2020. 6

[55] Guanshuo Wang, Yufeng Yuan, Xiong Chen, Jiwei Li, and Xi Zhou. Learning discriminative features with multiple granularities for person re-identification. In *ACMMM*, pages 274–282, 2018. 6

[56] Yicheng Wang, Zhenzhong Chen, Feng Wu, and Gang Wang. Person re-identification with cascaded pairwise convolutions. In *CVPR*, pages 1470–1478, 2018. 2

[57] Zhikang Wang, Lihuo He, Xiaoguang Tu, Jian Zhao, Xinbo Gao, Shengmei Shen, and Jiashi Feng. Robust video-based person re-identification by hierarchical mining. *T-CSVT*, 2021. 1

[58] Zhongdao Wang, Luming Tang, Xihui Liu, Zhuliang Yao, Shuai Yi, Jing Shao, Junjie Yan, Shengjin Wang, Hongsheng Li, and Xiaogang Wang. Orientation invariant feature embedding and spatial temporal regularization for vehicle re-identification. In *ICCV*, pages 379–387, 2017. 1

[59] Zhongdao Wang, Jingwei Zhang, Liang Zheng, Yixuan Liu, Yifan Sun, Yali Li, and Shengjin Wang. Cycas: Self-supervised cycle association for learning re-identifiable descriptions. In *ECCV*, pages 72–88, 2020. 6

[60] Longhui Wei, Shiliang Zhang, Wen Gao, and Qi Tian. Person transfer gan to bridge domain gap for person re-identification. In *CVPR*, pages 79–88, 2018. 2, 5, 12

[61] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *ECCV*, pages 3–19, 2018. 4

[62] Yuhui Xu, Lingxi Xie, Xiaopeng Zhang, Xin Chen, Guo-Jun Qi, Qi Tian, and Hongkai Xiong. Pc-darts: Partial channel connections for memory-efficient architecture search. In *ICLR*, 2019. 2

[63] Cheng Yan, Guansong Pang, Xiao Bai, Changhong Liu, Ning Xin, Lin Gu, and Jun Zhou. Beyond triplet loss: person re-identification with fine-grained difference-aware pairwise loss. *T-MM*, 2021. 1

[64] Mang Ye, Jianbing Shen, Gaojie Lin, Tao Xiang, Ling Shao, and Steven C.H. Hoi. Deep learning for person re-identification: A survey and outlook. *T-PAMI*, 2021. 1

[65] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z. Li. Deep metric learning for person re-identification. In *ICPR*, pages 34–39, 2014. 1

[66] Anguo Zhang, Yueming Gao, Yuzhen Niu, Wenxi Liu, and Yongcheng Zhou. Coarse-to-fine person re-identification with auxiliary-domain classification and second-order information bottleneck. In *CVPR*, pages 598–607, 2021. 6

[67] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, Xin Jin, and Zhibo Chen. Relation-aware global attention for person re-identification. In *CVPR*, pages 3186–3195, 2020. 2, 6, 7

[68] Aihua Zheng, Xianmin Lin, Jiacheng Dong, Wenzhong Wang, Jin Tang, and Bin Luo. Multi-scale attention vehicle re-identification. *Neural Computing and Applications*, 32(23):17489–17503, 2020. 1, 3

[69] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *ICCV*, pages 1116–1124, 2015. 2, 5, 13

[70] Liang Zheng, Yi Yang, and Alexander G Hauptmann. Person re-identification: Past, present and future. *arXiv preprint arXiv:1610.02984*, 2016. 1

[71] Liang Zheng, Yi Yang, and Qi Tian. Sift meets cnn: A decade survey of instance retrieval. *T-PAMI*, 40(5):1224–1244, 2017. 1

[72] Meng Zheng, Srikrishna Karanam, Ziyan Wu, and Richard J Radke. Re-identification with consistent attentive siamese networks. In *CVPR*, pages 5735–5744, 2019. 3, 14

[73] Zhedong Zheng, Xiaodong Yang, Zhiding Yu, Liang Zheng, Yi Yang, and Jan Kautz. Joint discriminative and generative learning for person re-identification. In *CVPR*, pages 2138–2147, 2019. 6

[74] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *AAAI*, pages 13001–13008, 2020. 6

[75] Kaiyang Zhou, Yongxin Yang, Andrea Cavallaro, and Tao Xiang. Omni-scale feature learning for person re-identification. In *ICCV*, pages 3702–3712, 2019. 1, 2, 3, 4, 5, 6

[76] Sanping Zhou, Fei Wang, Zeyi Huang, and Jinjun Wang. Discriminative feature learning with consistent attention regularization for person re-identification. In *ICCV*, pages 8040–8049, 2019. 3, 14

[77] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In *CVPR*, pages 8697–8710, 2018. 2

[78] Yang Zou, Xiaodong Yang, Zhiding Yu, BVK Vijaya Kumar, and Jan Kautz. Joint disentangling and adaptation for cross-domain person re-identification. In *ECCV*, pages 87–104, 2020. 1

Table 7. The detailed interaction operation comparison between MSINet, MSINet-VR and MSINet-S. N: None; E: Exchange; G: Channel Gate; A: Cross Attention.

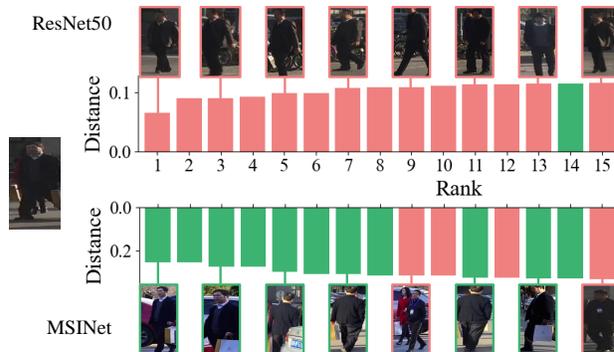| Model | #1 | | #2 | | #3 | | #4 | | #5 | | #6 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MSINet | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| | G | G | E | G | A | G | G | N | G | A | E | A |
| MSINet-VR | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| | G | G | E | A | G | A | G | A | A | A | E | A |
| MSINet-S | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| | E | A | G | A | A | A | G | A | E | A | E | A |



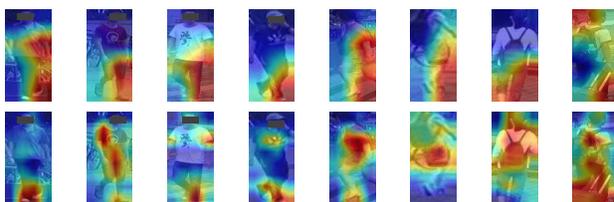Figure 7. Example top-20 retrieved sequences comparison on MSMT17. Best viewed in color.



Figure 8. Example activation maps of ResNet50 (the first row) and our proposed MSINet (the second row) trained on Market-1501 dataset. Best viewed in color.

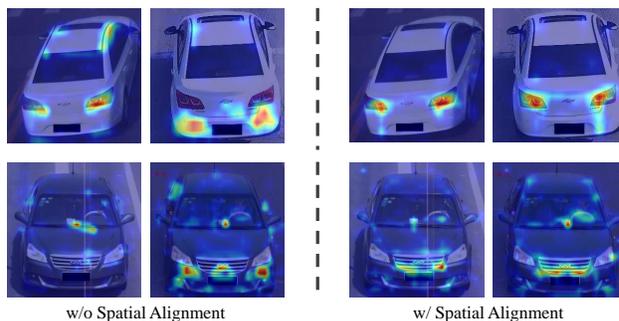

w/o Spatial Alignment       w/ Spatial Alignment

Figure 9. Example activation maps of MSINet trained on the task of VR→VID. With the Spatial Alignment Module, the model is capable to focus consistently on specific areas confronted with images from different sources. Best viewed in color.

# 7. Search on VeRi-776

We select a training-validation ratio of 60%-80% in the searching process on MSMT17 [60]. Without changing any specific configurations, we directly search for the rational interaction operations on VeRi-776 [38, 39] dataset. The searched architecture is denoted as MSINet-VR. We compare the structure of MSINet and MSINet-VR in Tab 7. Generally, the two searched architecture have common characteristics: Channel Gate is preferred in shallow layers, while Cross Attention is employed for more thorough information interaction in deep layers.

Quantitatively, we also conduct relevant supervision and cross-domain experiments with MSINet-VR in Tab. 8. All the experiment configurations are kept the same as those of MSINet training. Although there are some fluctuations, generally MSINet-VR has similar performance to MSINet, and the retrieval accuracy still surpasses ResNet50 [20, 42] by a large margin.

# 8. Search with Different Overlap Ratios

With the identities of training and validation sets unbound, we conduct a series of experiments utilizing different data separation ratios in Tab. 9 to find the appropriate interaction operations for the network. Firstly, we separate the training and validation sets completely with no identity overlaps. It can be observed that a balanced train-validation ratio generally brings better performance. For the two extremes of data distribution, an over-small validation set makes the architecture optimizer stuck in local minima and achieves poor performance. On the contrary, an over-large validation set brings no severe damage to the architecture search process, despite that the model is still not optimal. It demonstrates that abundant validation data is essential for ReID NAS.

Secondly, we randomly select part of identities, and evenly divide their images into the training and validation sets. The experiment results suggest that having a relatively small proportion of overlapped identities, whose images have been partly utilized for model parameter update, stabilizes the searching process and leads to a better architecture. However, when the overlap increases to a certain extent, the resemblance between the training and validation sets will bring negative influence to the ReID architecture search. As a comparison, we conduct the search task with traditional NAS scheme where a linear classification layer and cross entropy loss are employed for the training and validation data, the searched model of which performs worse than our proposed TCM.

Combined with above rules and the model performance, we select the architecture searched with the train-validation split of 60%-80% as the proposed MSINet.

Table 8. Supervised performance on object ReID datsets. The results in the top part are trained from scratch, and those in the bottom part are pre-trained on ImageNet in advance.

| Method | Params | Inference Time | M | | MS | | VR | | VID | | MS→M | | VR→VID | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | R-1↑ | mAP↑ | R-1↑ | mAP↑ | R-1↑ | mAP↑ | R-1↑ | R-5↑ | R-1↑ | mAP↑ | R-1↑ | R-5↑ |
| ResNet50* [42] ∼ | 24M | 1x | 85.7 | 68.3 | 48.0 | 25.7 | 92.8 | 69.9 | 70.6 | 76.6 | - | - | | |
| MSINet | 2.3M | 0.71x | 94.6 | 87.0 | **76.0** | **52.5** | **95.9** | 75.0 | **76.5** | **89.8** | - | - | - | - |
| MSINet-VR | 2.3M | 0.71x | **94.7** | **87.4** | 75.5 | 51.8 | **95.9** | **76.0** | 75.2 | 86.3 | - | - | - | - |
| ResNet50* [42] | ∼24M | 1x | 94.5 | 85.9 | 75.5 | 50.4 | 94.5 | 73.6 | 76.5 | 89.9 | 58.8 | 31.8 | 42.8 | 61.9 |
| MSINet | 2.3M | 0.71x | 95.3 | **89.6** | **81.0** | **59.6** | 96.8 | **78.8** | **77.9** | **91.7** | **74.9** | **46.2** | 48.0 | 65.6 |
| MSINet-VR | 2.3M | 0.71x | **95.4** | 89.0 | 80.1 | 57.5 | **97.0** | 78.6 | 77.3 | 91.3 | 72.9 | 44.8 | **48.5** | **66.2** |

Table 9. The evaluation results of models searched with different training-validation identity ratios on the MS dataset. * indicates that the model is searched with the softmax loss.

| w/o overlap | | | | w/ overlap | | | |
|---|---|---|---|---|---|---|---|
| | | MS | | | | MS | |
| train (%) | valid (%) | R-1↑ | mAP↑ | train (%) | valid (%) | R-1↑ | mAP↑ |
| 90 | 10 | 70.6 | 46.2 | 60 | 60 | 75.3 | 51.2 |
| 75 | 25 | 71.4 | 46.9 | 40 | 80 | 75.1 | 50.5 |
| 67 | 33 | 75.5 | 51.5 | 60 | 80 | **76.0** | **52.5** |
| 50 | 50 | 74.7 | 50.4 | 80 | 60 | 74.8 | 51.0 |
| 33 | 67 | 74.7 | 50.8 | 80 | 80 | 74.4 | 50.3 |
| 25 | 75 | 74.9 | 50.6 | 100 | 100 | 74.4 | 50.1 |
| 10 | 90 | 75.4 | 50.7 | 100* | 100* | 74.0 | 49.8 |

## 9. Search with Softmax Loss

We further compare the detailed interaction operations between MSINet and the architecture searched under traditional NAS scheme, where softmax loss and a unified linear classification layer are utilized for the training and validation sets [36] (denoted as MSINet-S) in Tab. 7. Compared with MSINet and MSINet-VR, where direct information exchange mainly appears at deep layers, MSINet-S contains a large amount of Exchange and Cross Attention along the whole network. The over-frequent information exchange fails to focus on discriminative features. It also validates the effectiveness of our proposed Twins Contrastive Mechanism on searching for architectures suitable for ReID.

## 10. Visualization Results

Some additional visualization results are illustrated to further manifest the effectiveness of our proposed architecture. Firstly, we visualize an example comparison of the top-20 retrieved sequences between ResNet50 and MSINet on MSMT17 in Fig. 7. ResNet50 mainly focus on general appearance information, while our proposed MSINet concentrates on discriminative distinctions, the hand bag in this case. Even though positive samples have large appearance differences from the query image, MSINet is still capable to distinguish them.

Secondly, example activation maps of ResNet50 and our proposed MSINet on Market-1501 [69] are visualized in Fig. 8. ResNet50 mainly focuses on the right part of the image, including some background areas. Our proposed MSINet, oppositely, is capable to dynamically focus on discriminative distinctions of each image.

Thirdly, to intuitively demonstrate the effectiveness of the Spatial Alignment Module (SAM) on enhancing the attention consistency of the model confronted with images from different sources, we visualize example activation maps on the task of VR→VID. As shown in Fig. 8, without alignment, the model can have different activated positions on different images of the same identity, even if they share similar appearances.

## 11. Comparison and advantages to OSNet

(1) OSNet simply sums up the features of each branch, without detailed exploration on the interaction between branches. In comparison, MSINet practically select rational interaction operations for different network layers. Consequently, MSINet surpasses OSNet not only in supervised, but also in domain generalization performance by a large margin. (2) OSNet contains 4 branches with different receptive field scales, where there exists certain parameter redundancy. We validated in the early exploring that removing

| Method | M→MS | | VID→VR | |
|---|---|---|---|---|
| | R-1↑ | mAP↑ | R-1↑ | mAP↑ |
| OSNet | 21.2 | 7.2 | 69.2 | 32.0 |
| MSINet | **22.4** | **8.3** | **72.1** | **33.8** |

Table 10. Additional Cross-domain Experiments

| Method | M | | VR→VID | |
|---|---|---|---|---|
| | R-1↑ | mAP↑ | R-1↑ | R-5↑ |
| SAM | **95.5** | **89.9** | **49.0** | **66.8** |
| SAM-softmax | 95.0 | 89.1 | 48.2 | 65.9 |

Table 11. Ablation study on softmax operation.

the branches with receptive field scales of 3 and 5 has little influence to the model performance. MSINet reduces the number of branches, and increases the scale difference between two branches, which increases the parameter amount by a little bit but significantly reduces the inference time.

## 12. Detailed Analysis on SAM

We compare the proposed SAM module to some previous attention-based methods and analyze it in detail. [76] regularizes the attention generated at different network layers for the same image; [72] explicitly enforces the longitudinal activation distribution to be the same for two images, which may lead to misalignment if the objects are not properly detected. For negative samples, there can be many different hints for recognition, some of which might be inappropriate, such as the backgrounds. By aggregating the information from different negative samples, the network is driven to only focus on discriminative regions. The motivation of SAM is different from the above two works. For the in-domain setting, the camera condition diffrences are directly addressed by supervised learning. Thus, SAM brings limited improvements, yet doesn't defect the performance, compared to techniques like instance normalization.

## 13. More Ablation Study

**Additional cross-domain experiments.** We add the M→MS and VID→VR experiment results to Tab. 10. MSINet surpasses OSNet on all metrics.

**Ablation study on softmax operation.** SAM aligns the activation values in the feature map, where the discriminative positions are actually matched between different samples. As the "Mutual Conv" operation is non-parametric, it is not proper to apply the softmax-squeezed position attention values for direct alignment, which may result in scale inconsistency. The experiment results in Tab. 11 also suggest slight influence on this detail.

## 14. Limitations and Future Work

The designed interaction operations only include forward and exchange in the direct and attention forms, which restricts the size of the search space. In the future works, there are still exploration room for more elaborate and complicated interaction operations and search spaces. There are still exploration room for more complicated search schemes and spaces.