

Neural Feature Search for RGB-Infrared Person Re-Identification

Yehansen Chen^{1*}, Lin Wan^{1*}, Zhihang Li^{2†}, Qianyan Jing¹, Zongyuan Sun¹

¹School of Geography and Information Engineering, China University of Geosciences, Wuhan, China

²School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing, China

Abstract

RGB-Infrared person re-identification (RGB-IR ReID) is a challenging cross-modality retrieval problem, which aims at matching the person-of-interest over visible and infrared camera views. Most existing works achieve performance gains through manually-designed feature selection modules, which often require significant domain knowledge and rich experience. In this paper, we study a general paradigm, termed Neural Feature Search (NFS), to automate the process of feature selection. Specifically, NFS combines a dual-level feature search space and a differentiable search strategy to jointly select identity-related cues in coarse-grained channels and fine-grained spatial pixels. This combination allows NFS to adaptively filter background noises and concentrate on informative parts of human bodies in a data-driven manner. Moreover, a cross-modality contrastive optimization scheme further guides NFS to search features that can minimize modality discrepancy whilst maximizing inter-class distance. Extensive experiments on mainstream benchmarks demonstrate that our method outperforms state-of-the-arts, especially achieving better performance on the RegDB dataset with significant improvement of 11.20% and 8.64% in Rank-1 and mAP, respectively.

1. Introduction

Person re-identification (ReID) aims to match the person-of-interest over non-overlapping camera views [54, 67, 44, 73, 60], serving as a central part of intelligent video surveillance systems. Currently, most conventional ReID methods concentrate efforts on visible images-based cross-view matching tasks [29, 18, 46, 43, 36], which cannot adapt well to illumination variations in real-world scenarios (e.g., low lighting environments at nighttime). Motivated by this challenge, associating RGB and infrared (IR) pedestrian images captured by dual-mode cameras for cross-modality image retrieval, a.k.a. RGB-IR ReID, has drawn much attention in vision community [58, 64, 51, 52, 66].

*Equally-contributed first authors

†Corresponding author

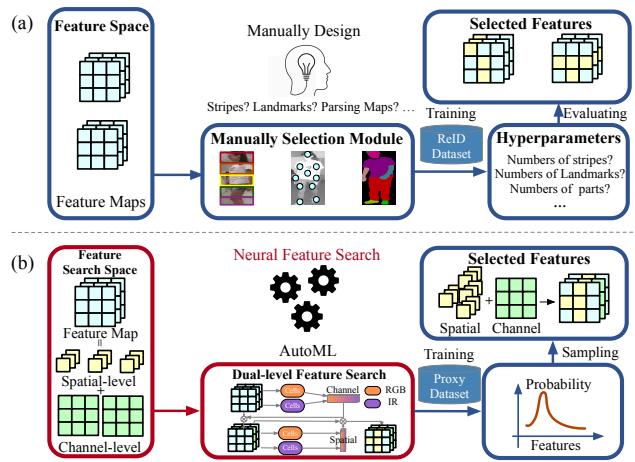


Figure 1. Comparison of hand-crafted and automated feature selection strategies. (a) Manually designing task-specific modules to select identity-related features. (b) NFS automatically derives the optimal feature subset from a dual-level feature search space.

Due to intrinsically different imaging mechanisms, RGB-IR ReID suffers from undesired visual discrepancy between visible and infrared images, which makes appearance cues such as colors and textures unreliable or even misleading for the matching task [32, 52, 58]. Moreover, such modality divergence also exacerbates the already large intra-class variations caused by diverse camera viewpoints, person poses, partial occlusions, and background clutter [57, 14, 22], making it even harder to align images of the same identity. In an effort to minimize the modality gap, cross-modality image synthesis methods [51, 52, 22] typically leverage generative adversarial networks (GANs) to transfer stylistic properties between modalities to synthesize fake RGB/IR images. But it is non-trivial to preserve identity information for generated RGB images due to lack of color information in their IR counterparts [11]. Another line of shared feature learning approaches [64, 69, 66, 63] utilize convolutional neural networks (CNNs) to perform cross-modality feature alignment. One representative model-of-choice is the two-stream network [63, 68, 56], which includes modality-specific shallow layers and shared deeper layers to learn a common feature

space [58]. On the strength of two-stream structures, several studies exploit ReID discriminative constraints, e.g., triplet loss [67, 65, 63] or ranking loss [69, 14], to supervise the network to mine identity-related cues. They are all committed to learning a better distance metric that enhances the performance of similarity-based retrieval and have achieved significant success in recent years [66].

To our understanding, whether image synthesis approaches or shared feature learning techniques, the crux of ReID solutions is always to find sufficiently high-quality discriminative features for matching and retrieval. To achieve this goal, state-of-the-art methods typically introduce partition stripes [66], human landmarks [49], parsing maps [19], and body contour sketches [62] to discourage irrelevant features whilst preserving the identity-related ones (Fig. 1(a)). However, it is really tough and time-consuming to manually design a *one-fit-all* feature selection module against all sorts of intra- and inter-modality variations, leading to unsatisfactory performance of human-guided feature selection mechanisms. Driven by the above observations, a question arises: *Is there a data-driven feature selection manner without much requirement for human interference?* Recent advances on automated machine learning (AutoML) [17] may provide a positive answer. Using Neural Architecture Search (NAS) [72], Quan *et al.* [41] automatically generate fast and effective CNNs whose performance is on par with hand-crafted architectures in single-modality ReID tasks. This progress inspires the idea of neuron-powered automatic feature selection discussed in this paper.

With the idea in mind, we investigate RGB-IR ReID from a *one-fit-all* feature search perspective. Different from existing manually-designed [66, 32, 4] or NAS-generated network structures [41], our goal is to pursue better ReID performance by discovering discriminative features with data-driven search neurons. To this end, we cast feature selection as a bilevel optimization problem [28] (i.e., deriving the optimal feature subset from the best feature learning results) and propose a novel paradigm, *Neural Feature Search* (NFS, Fig. 1(b)). Starting from the hierarchical feature extraction mechanism of CNNs [39], NFS includes a dual-level feature search space where each feature map is decomposed in terms of pixel and channel dimensions. This allows feature selection operations to be jointly performed in a mutually reinforcing manner—channel-level search identifies relevant response maps from the global view while pixel-level search scans every spatial position to selectively process local part features of a person. To improve the search efficiency, we utilize reparameterization tricks [31] to relax the search space to be continuous, which makes the optimization of search neurons compatible with stochastic gradient descent (SGD). Considering the inherent modality discrepancy issue of RGB-IR ReID, a cross-modality contrastive optimization scheme is further introduced as a su-

pervision signal that discourages irrelevant features whilst encouraging modality-invariant cues during the searching process. Extensive experiments show that NFS significantly outperforms the state-of-the-arts by 12.01% and 11.20% gains of Rank-1 accuracy on SYSU-MM01 (*multi-shot, all search mode*) and RegDB (*visible-to-infrared mode*) benchmarks, respectively. To summarize, this paper brings three main contributions:

- We propose an AutoML-powered neural feature search method for RGB-IR person re-identification, which automates the feature selection process to substantially improve the matching accuracy with less human efforts. To our best knowledge, this is one of the first attempts to utilize automatic feature selection techniques for cross-modality ReID.
- We introduce a novel feature search space allowing both spatial and channel-wise feature selection. Based on this search space, we present an efficient feature search algorithm embedded with a cross-modality contrastive optimization mechanism, effectively tackling the modality discrepancy in RGB-IR ReID.
- Extensive experiments on two mainstream RGB-IR ReID benchmarks demonstrate the superiority of NFS compared with previous state-of-the-arts.

2. Related Work

RGB-based Person ReID. RGB-based person ReID studies mainly focus on handling intra-class variations of pose [29], scale [18], and background clutter [46] presented in visible images. Nowadays, substantial research efforts [18, 50, 30, 8, 61] have been devoted to deep learning-based ReID for more effective feature learning and alignment. For example, graph neural networks [43, 36] make full use of relationships between global embedding vectors to map images into a discriminative feature space. Self-attention based methods [48, 33] explore pixel similarities to let the network concentrate on informative biometrics such as face against the background clutter. Apart from global feature representation learning, several local feature learning approaches [49, 71] also employ pretrained pose estimation models to decompose the human body into landmarks or parsing maps and perform fine-grained feature alignment over pose changes and occlusions. Although having achieved considerable success in reducing intra-class variations, most existing single-modality ReID methods are ill-suited for cross-modality image retrieval in poor lighting environments [68, 65, 63].

RGB-Infrared Person ReID. In addition to intra-class variations, RGB-IR ReID also considers the modality discrepancy issue caused by different wavelength ranges of visible and infrared cameras [52]. Current RGB-IR ReID

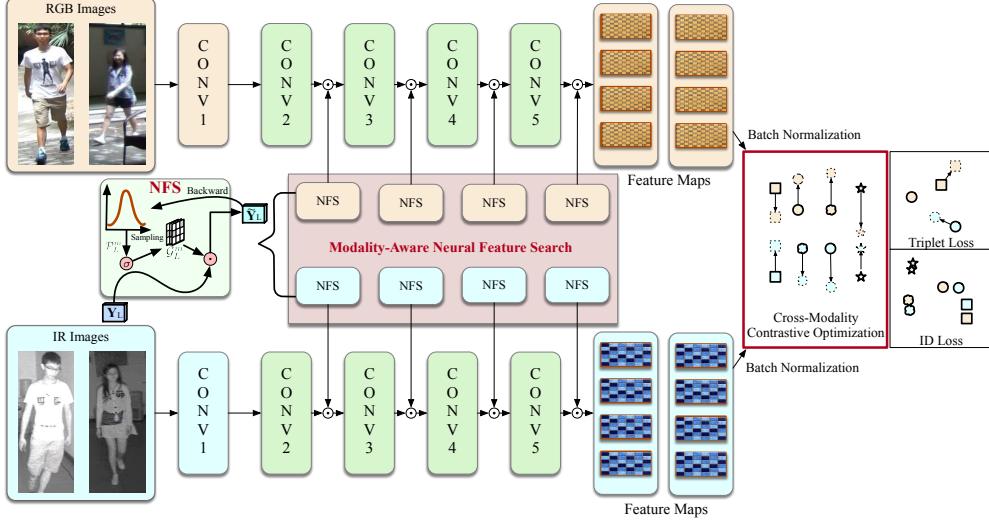


Figure 2. An overview of our NFS paradigm. It combines modality-aware search cells and cross-modality contrastive optimization mechanism to conduct automated feature selection on two-stream CNNs based feature space. Each learnable search cell is jointly optimized with network parameters to derive the optimal feature subset in every shared block. The cross-modality contrastive optimization mechanism further enables NFS to search modality-invariant features that can minimize modality discrepancy while maximizing inter-class distance.

researches mainly resort to either GAN-based [51, 52, 5, 4, 70, 55] or shared feature learning approaches [69, 64, 66, 58, 32] to handle both intra- and inter-modality variations.

For GAN-based approaches, early attempts [5] usually adopt adversarial training strategies to reduce the distribution divergence of modality-specific features. Along a somewhat different line, Wang *et al.* [51] leverage GANs to transfer stylistic properties of IR images to their RGB counterparts for jointly pixel and feature alignment. Several studies also apply pair-wise pixel alignment [52], feature disentanglement [4], or intermediate modality generation [22, 55] to further eliminate appearance differences across modalities. However, it is non-trivial to accurately choose the suitable target for style transfer [32], which may lead to identity inconsistency during the complicated adversarial training process [11]. As for the shared feature learning category, Wu *et al.* [58] first contribute a large benchmark dataset (SYSU-MM01) and propose a one-stream zero-padding network for RGB-IR image matching. Nowadays, two-stream CNNs based methods are dominating the cross-modality person ReID community. For instance, some recent studies extend two-stream CNNs with deep metric learning [64, 69, 14, 63, 65, 68, 45] or the attention mechanism [66, 56] to learn modality-sharable representations against both modality discrepancy and high sample noises. Several works [10, 32] also employ modality-specific or modality-aware learning avenues to perform cross-modality identity recognition at the classifier level.

Neural Architecture Search. Recent years have witnessed a growing body of Neural Architecture Search

(NAS) researches [9, 23] that have achieved considerable success in various domains, e.g., image classification [42], semantic segmentation [3], object detection [25], ReID [41], and multi-modal fusion [39], etc. Generally, NAS aims to automatically search optimal operations or topology of deep neural networks for specific learning tasks. They first construct a task-oriented search space that defines which architectures can be discovered in principle. Based on the search space, different search strategies, including reinforcement learning-based methods [12, 40], evolutionary methods [27, 42], gradient-based ones [28, 24], and Monte Carlo Tree Search (MCTS) approaches [37, 53], are proposed and prove effective for improving both sample efficiency and model performance. Inspired by the basic idea of NAS, we present a *one-fit-all* feature selection strategy for RGB-IR person ReID, revolutionizing manually-crafted feature selection components in the existing literature.

3. Methodology

Fig. 2 presents an overview of our proposed method. On the basis of a two-stream network (Section 3.1), NFS mainly includes a dual-level search space for spatial and channel-wise feature selection, and a differentiable feature search algorithm (Section 3.2) governed by the cross-modality contrastive optimization scheme (Section 3.3) to prune discriminative cues fast and accurately.

3.1. Baseline RGB-IR Person ReID

We adopt the two-stream CNN employed in [69, 67, 66] as the baseline network[†]. To capture modality-invariant information, parameters of the first convolutional block are independent for each modality, while the other layers are shared to learn discriminative features [58]. After the last convolutional layer with global average pooling, a shared batch normalization layer is used to attain final representations for heterogeneous images. During the training phase, we aim to minimize the following baseline loss function:

$$\mathcal{L}_b = \mathcal{L}_{id} + \mathcal{L}_{tri}, \quad (1)$$

where \mathcal{L}_{id} is the softmax cross-entropy loss and \mathcal{L}_{tri} is the weighted regularization triplet (WRT) loss [67].

3.2. Modality-aware Neural Feature Search

Unlike NAS approaches searching optimal topology and operations for a top-performing architecture [28, 59, 40], NFS searches for identity-related features from a CNN-based feature space. In this paper, we cast the automatic feature search as a hyperparameter learning task, where the search hyperparameters and network weights are jointly optimized to derive the optimal discriminative feature subset. It can be formulated as a bilevel optimization problem [28]:

$$\min_{y \in Y} \min_W \mathcal{L}(y, W). \quad (2)$$

Given output feature maps Y , we seek to discover a subset of features $y \in Y$, which can minimize the loss $\mathcal{L}(y, W)$ after optimizing network weights W . Here, we highlight two key points of solving the problem: a dual-level feature search space and an efficient search algorithm.

Dual-level Feature Search Space. The search space defines what neural architectures might be discovered in principle [28, 59], which plays a crucial role in high-performance NAS. Similarly, NFS is also closely related to a well-designed feature search space that covers as many as possible identity-related cues. As discriminative features are mainly extracted by the shared part of the baseline model [58], we establish a search space including all feature candidates extracted by every shared convolutional block. More formally, a shared block L takes a feature map $X \in \mathbb{R}^{C_{in} \times W \times H}$ as input and outputs another fine-grained feature map $Y_L \in \mathbb{R}^{C_{out} \times \frac{W}{2} \times \frac{H}{2}}$, i.e.,

$$Y_L(p) = \sum_{p' \in R_k} W_c(p') X(p + p') \quad p \in \Omega, \quad (3)$$

where C , W , and H represent the number of channels, width, and height, respectively. p denotes a specific pixel position, R_k is the support region of kernel with size k ,

[†]<https://github.com/mangye16/Cross-Modal-Re-ID-baseline>

$W_c \in \mathbb{R}^{C_{in} \times C_{out} \times k \times k}$ represents convolution weights, $\Omega = \{(i, j) | i \leq W, j \leq H, i, j \in \mathbb{Z}^+\}$ is the spatial domain of Y_L . The union of Y_L forms the vanilla feature space $Y = \{Y_L | L \in \{1, \dots, N\}\}$, and N is the number of blocks.

Motivated by the fact that discriminative features present modality-specific distributions in spatial and channel dimensions of Y_L [4], we introduce modality-aware search cells to decompose the original feature space Y into two subspaces: pixel-level subspace and depth-level subspace. The former includes vectors in each spatial position that describe local patches of an input image. The latter contains multiple detector response maps that globally reflect particular properties of the image content. Fig. 2(Left) illustrates how the search cell exactly works. Given the output feature map Y_L , we first initialize a set of parameters \mathcal{P}_L^m with a uniform distribution to map features from modality m into a specific probability field. As for pixel-level feature search, $\mathcal{P}_L^m \in \mathbb{R}^{C_{out} \times \frac{W}{2} \times \frac{H}{2}}$ covers every pixel in the spatial domain. And for depth-level search, $\mathcal{P}_L^m \in \mathbb{R}^{C_{out}}$ contains all channels of Y_L . During the searching process, the probability field is activated by a sigmoid function, denoted as $\tilde{\mathcal{P}}_L^m = \sigma(\mathcal{P}_L^m)$, to indicate the possibilities of features at corresponding positions are informative to distinguish different persons. Then, a binary search gate $\mathcal{G}_L^m(p)$ is generated based on $\tilde{\mathcal{P}}_L^m$ to determine whether the pixel at position p should be selected. By passing Y_L through all search gates, the output activation map \tilde{Y}_L can be formulated as:

$$\tilde{Y}_L(p) = \begin{cases} Y_L(p), & \mathcal{G}_L^m(p) = 1 \\ 0, & \mathcal{G}_L^m(p) = 0. \end{cases} \quad (4)$$

Here, Eq. 2 is transformed into an optimization problem with all search gates \mathcal{G}^m as the upper-level variables and the network weights W as lower-level ones, that is:

$$\begin{aligned} \min_{\mathcal{G}^m} \quad & \mathcal{L}_{val}(W^*(\mathcal{G}^m), \mathcal{G}^m) \\ \text{s.t.} \quad & W^*(\mathcal{G}^m) = \arg \min_W \mathcal{L}_{train}(W, \mathcal{G}^m), \end{aligned} \quad (5)$$

where \mathcal{L}_{train} and \mathcal{L}_{val} denote the training loss and the validation loss (Eq. 14), respectively.

Search Algorithm. As the search space is discrete and large-scale, finding the optimal feature set through brute-force enumeration is much inefficient. To tackle this obstacle, we utilize reparameterization tricks [28] to relax the search space to be continuous and directly improve the search efficiency via SGD. We assume that feature selection is essentially a binary classification problem and thus exploiting a continuous Bernoulli distribution [31] to simulate stochastic discrete sampling with $\tilde{\mathcal{P}}^m \in (0, 1)$. Based on \mathcal{G}^m , the sampled features \mathcal{X}^m are as:

$$\begin{aligned} \mathcal{X}^m \sim Bernoulli(\tilde{\mathcal{P}}^m) \iff & p(x^m | \tilde{\mathcal{P}}^m) \propto \hat{p}(x^m | \tilde{\mathcal{P}}^m) \\ = & (\tilde{\mathcal{P}}^m)^{x^m} (1 - \tilde{\mathcal{P}}^m)^{1-x^m}, \end{aligned} \quad (6)$$

$$p(x^m | \tilde{\mathcal{P}}^m) = C(\tilde{\mathcal{P}}^m)(\tilde{\mathcal{P}}^m)^{x^m}(1 - \tilde{\mathcal{P}}^m)^{1-x^m}, \quad (7)$$

$$C(\tilde{\mathcal{P}}^m) = \begin{cases} \frac{2\tanh^{-1}(1-2\tilde{\mathcal{P}}^m)}{1-2\tilde{\mathcal{P}}^m}, & \text{if } \tilde{\mathcal{P}}^m \neq 0.5 \\ 2, & \text{otherwise.} \end{cases} \quad (8)$$

After relaxation, \mathcal{G}^m and W can be jointly optimized using the straight-through estimator (STE) [1]:

$$\nabla_{\mathcal{G}^m} \mathcal{L}_{val} \approx \nabla_{\tilde{\mathcal{P}}^m} \mathcal{L}_{val}. \quad (9)$$

Finally, Eq. 5 is solved with Neural Feature Search outlined in Alg. 1. We first search for identity-related features by iteratively optimizing \mathcal{G}^m with \mathcal{L}_{val} and W with \mathcal{L}_{train} . After obtaining the optimal search gates, the whole network is trained and evaluated following the standard feature learning paradigm of RGB-IR ReID [66].

Algorithm 1 NFS - Neural Feature Search

Input: the search parameters $\tilde{\mathcal{P}}^m$ and the network weights W ; the training set \mathbb{D}_T and the testing set \mathbb{D}_E

Output: the trained network and the optimal feature set

- 1: Randomly split \mathbb{D}_T into the search training set \mathbb{D}_{train} and the search validation set \mathbb{D}_{val}
 - 2: **while** not converged **do**
 - 3: Update W by descending $\nabla_W \mathcal{L}_{train}(W, \tilde{\mathcal{P}}^m)$ on \mathbb{D}_{train}
 - 4: Update search gates \mathcal{G}^m by descending $\nabla_{\tilde{\mathcal{P}}^m} \mathcal{L}_{val}(W - \nabla_W \mathcal{L}_{train}(W, \tilde{\mathcal{P}}^m), \tilde{\mathcal{P}}^m)$ on \mathbb{D}_{val}
 - 5: Derive optimal \mathcal{G}^m based on $\tilde{\mathcal{P}}^m$
 - 6: Train the network weight W with the derived \mathcal{G}^m on \mathbb{D}_T
 - 7: Evaluate the network and \mathcal{G}^m on \mathbb{D}_E
-

3.3. Cross-Modality Contrastive Optimization

Apart from the search efficiency, how to supervise search cells to select more informative features is also important for NFS. Unlike close-set classification tasks, RGB-IR ReID is an open-set problem where identities in testing are different from those in training. In such a scenario, the selected features of ‘seen’ and ‘unseen’ classes may be tangled in the feature space. Meanwhile, the appearance discrepancy between RGB and IR images often enlarges the feature distribution variance of each class, leading to fuzzy decision boundaries in identity recognition problems.

Here, we attend to decrease the feature distribution variance from an invariant feature selection perspective. To this end, we introduce a ReID-oriented optimization criterion that can eliminate modality discrepancy and maximize the inter-class distance simultaneously. The basic idea comes from recent advances on contrastive learning [47, 21, 15], which aim to attract positive pairs whilst repulsing negative ones [2]. Given a training batch $B = \{(i_{rgb}, i_{ir}) | i_{rgb} \in \mathcal{I}_{rgb}, i_{ir} \in \mathcal{I}_{ir}\}$, the half of which are RGB images \mathcal{I}_{rgb} while the others are their IR counterparts \mathcal{I}_{ir} , we randomly

arrange their embedding vectors \vec{i}_{rgb} and \vec{i}_{ir} into multiple cross-modality pairs $(\vec{i}_{rgb}, \vec{i}_{ir})$ and generate pair-wise labels according to their identities $ID(\vec{i}_{rgb})$ and $ID(\vec{i}_{ir})$:

$$Label = \begin{cases} 1, & ID(\vec{i}_{rgb}) = ID(\vec{i}_{ir}) \\ 0, & ID(\vec{i}_{rgb}) \neq ID(\vec{i}_{ir}). \end{cases} \quad (10)$$

For each positive image pair, we seek to minimize the distance between them, so that the modality discrepancy and intra-class variations can be jointly eliminated. We evaluate the pair-wise distance in Euclidean space, which is widely applied in image retrieval [13, 67], i.e.,

$$D_E = \|\vec{i}_{rgb} - \vec{i}_{ir}\|_2. \quad (11)$$

On the contrary, for negative pairs, we aim to keep them far from each other for distinction. In order to make the optimization objective explicitly, we quantify the dissimilarity of each negative pair D_T with an explicit margin T :

$$D_T = \max(0, T - \|\vec{i}_{rgb} - \vec{i}_{ir}\|_2). \quad (12)$$

Taking all positive and negative pairs into account, the contrastive loss \mathcal{L}_c can be formulated as:

$$\mathcal{L}_c(Label, \vec{i}_{rgb}, \vec{i}_{ir}) = (Label)(D_E)^2 + (1 - Label)(D_T)^2. \quad (13)$$

The overall learning objective for NFS is a weighted summation of the baseline loss \mathcal{L}_b and cross-modality contrastive loss \mathcal{L}_c , defined as:

$$\mathcal{L} = \mathcal{L}_b + \lambda \mathcal{L}_c, \quad (14)$$

where λ is a trade-off coefficient to balance the influence of each learning objective.

4. Experiments

4.1. Datasets and Experimental Settings

Datasets. Our experiments are based on two standard real-world benchmarks for RGB-IR person ReID, named SYSU-MM01 [58] and RegDB [38], respectively. The SYSU-MM01 dataset contains images captured by four visible and two near infrared cameras in indoor and outdoor environments. Statistically, the training set includes 22,258 RGB and 11,909 IR images of 395 identities, and the query set involves 3,803 IR images of 96 identities. The gallery set has four versions according to different evaluation modes, including *all search* or *indoor search* and *single-shot* or *multi-shot*. Details of each mode can be found in [58]. The RegDB dataset contains 8,240 images of 412 identities, with 206 identities for training and the rest for testing. Each identity has 10 IR and 10 RGB images. We evaluate both *visible-to-infrared* and *infrared-to-visible* modes [51] by alternatively using all RGB/IR images as the gallery set.

Table 1. Comparison on the SYSU-MM01 dataset with Rank-1, 10, 20 (%) and mAP (%) evaluation metrics.

Method	All Search								Indoor Search							
	Single-shot				Multi-shot				Single-shot				Multi-shot			
	r1	r10	r20	mAP	r1	r10	r20	mAP	r1	r10	r20	mAP	r1	r10	r20	mAP
HOG [6]	2.76	18.3	31.9	4.24	3.82	22.8	37.6	2.16	3.22	24.7	44.5	7.25	4.75	29.2	49.4	3.51
LOMO [26]	3.64	23.2	37.3	4.53	4.70	28.2	43.1	2.28	5.75	34.4	54.9	10.2	7.36	40.4	60.3	5.64
Zero-Padding [58]	14.8	54.1	71.3	15.9	19.1	61.4	78.4	10.9	20.6	68.4	85.8	26.9	24.4	75.9	91.3	18.6
TONE+HCML [64]	14.3	53.2	69.2	16.2	-	-	-	-	-	-	-	-	-	-	-	-
BDTR [69]	17.0	55.4	72.0	19.7	-	-	-	-	-	-	-	-	-	-	-	-
D-HSME [14]	20.7	62.8	78.0	23.2	-	-	-	-	-	-	-	-	-	-	-	-
IPVT+MSR [20]	23.2	51.2	61.7	22.5	-	-	-	-	-	-	-	-	-	-	-	-
cmGAN [5]	27.0	67.5	80.6	27.8	31.5	72.7	85.0	22.3	31.6	77.2	89.2	42.2	37.0	80.9	92.1	32.8
D ² RL [55]	28.9	70.6	82.4	29.2	-	-	-	-	-	-	-	-	-	-	-	-
DGD+MSR [10]	37.4	83.4	93.3	38.1	43.9	86.9	95.7	30.5	39.6	89.3	97.7	50.9	46.6	93.6	98.8	40.1
JSIA-ReID [52]	38.1	80.7	89.9	36.9	45.1	85.7	93.8	29.5	43.8	86.2	94.2	52.9	52.7	91.1	96.4	42.7
AlignGAN [51]	42.4	85.0	93.7	40.7	51.5	89.4	95.7	33.9	45.9	87.6	94.4	54.3	57.1	92.7	97.4	45.3
AGW [67]	47.50	84.39	92.14	47.65	-	-	-	-	54.17	91.14	95.98	62.97	-	-	-	-
Xmodal [22]	49.92	89.79	95.96	50.73	-	-	-	-	-	-	-	-	-	-	-	-
DDAG [66]	54.75	90.39	95.81	53.02	-	-	-	-	61.02	94.06	98.41	67.98	-	-	-	-
NFS (Ours)	56.91	91.34	96.52	55.45	63.51	94.42	97.81	48.56	62.79	96.53	99.07	69.79	70.03	97.70	99.51	61.45

Table 2. Comparison on the RegDB dataset with Rank-1, 10, 20 (%) and mAP (%) evaluation metrics.

Method	Visible to Infrared				Infrared to Visible			
	r1	r10	r20	mAP	r1	r10	r20	mAP
Zero-Padding [58]	17.75	34.21	44.35	18.90	16.63	34.68	44.25	17.82
Tone + HCML [64]	24.44	47.53	56.78	20.88	21.70	45.02	55.58	22.24
BDTR [69]	33.56	58.61	67.43	32.76	32.92	58.46	68.43	31.96
D ² RL [55]	43.4	66.1	76.3	44.1	-	-	-	-
DGD+MSR [10]	48.43	70.32	79.95	48.67	-	-	-	-
JSIA-ReID [52]	48.1	-	-	48.9	48.5	-	-	49.3
D-HSME [14]	50.85	73.36	81.66	47.00	50.15	72.40	81.07	46.16
IPVT+MSR [20]	58.76	85.75	90.27	47.85	-	-	-	-
AlignGAN [51]	57.9	-	-	53.6	56.3	-	-	53.4
Xmodal [22]	62.21	83.13	91.72	60.18	-	-	-	-
DDAG [66]	69.34	86.19	91.49	63.46	68.06	85.15	90.31	61.80
NFS (Ours)	80.54	91.96	95.07	72.10	77.95	90.45	93.62	69.79

Evaluation Protocols. We follow standard evaluation protocols [58, 67] for RGB-IR ReID. Gallery and query images are from different modalities. The standard cumulated matching characteristics (CMC) curve and mean average precision (mAP) are used for performance evaluation.

Implementation Details. The proposed method is implemented in PyTorch and trained on an NVIDIA 2080Ti GPU. In order to facilitate comparisons with state-of-the-art ReID researches [67, 66, 32], we adopt the ResNet-50 [16] pretrained on ImageNet [7] as our backbone network. Following [66, 65, 63, 67], we set the stride of the last convolutional block as 1 for fine-grained feature maps. All images are resized to 288 × 144 then augmented with random cropping and horizontal flipping. We randomly sample 80% images from the original training set as the search training set and use the rest as the search validation set (Alg. 1, Line 1). We first make all search cells learnable to discover the optimal discriminative feature set. After obtaining the optimal feature set, we retrain the network on the original training set. At the training stage, we adopt a warm-up strategy [34] and optimize the two-stream CNN using SGD with 0.9 momentum during a total of 80 epochs. The initial learning rate is set to 0.1 and decays by 0.1 and 0.01 at the 16th and 50th epoch, respectively. Following [66], we randomly sample 8 identities with 4 RGB and 4 IR images per person, resulting in totally 64 images for each training batch.

4.2. Comparison with State-of-the-art Methods

In this subsection, we compare the proposed NFS with naive baselines as well as the state-of-the-art methods, including traditional feature extraction methods (HOG [6] and LOMO [26]); GAN-based models (cmGAN [5], D²RL [55], JSIA-ReID [52], AlignGAN [51], and Xmodal [22]); deep metric learning (BDTR [69], D-HSME [14], IPVT+MSR [20], and DGD+MSR [10]); and shared feature learning approaches (Zero-Padding [58], TONE+HCML [64], AGW [67], and DDAG [66]). Since most of them follow the standard evaluation protocols of the two experimental datasets, we directly use the original results from published papers for comparison.

Experimental results on SYSU-MM01 are shown in Table 1. We see that there is a significant performance decline when applying hand-crafted descriptors HOG and LOMO to cross-modality ReID, regardless of their promising capacities in general ReID tasks. Besides, image synthesis methods (AlignGAN, JSIA-ReID, Xmodal, and D²RL) perform better than traditional shared feature learning approaches (Zero-Padding and TONE+HCML), possibly owing to the effectiveness of pixel-level alignment. Specifically, recent methods such as AGW, DDAG, as well as our proposed NFS outperform typical GAN-based approaches. This is probably because it is ill-posed to transfer identity-

related information of IR images to generated RGB images. Notably, the proposed model achieves 56.91% Rank-1 identification rate and 55.45% mAP score in the most difficult *single-shot & all search* setting, which outperforms most of state-of-the-art methods by a large margin. Compared to DDAG based on the graph attention mechanism, NFS is much easier to implement and still presents better performance. Similar improvement can be observed in *multi-shot* modes. For example, our method largely surpasses AlignGAN with the improvement of 12.01% in Rank-1 and 14.66% in mAP, which demonstrates highly robustness when the gallery size increases.

Results on the RegDB dataset are listed in Table 2. Generally, performance of all methods is higher than that on SYSU-MM01, as images of RegDB present less intra-class variations [51]. Our approach greatly improves the state-of-the-art under both evaluation modes. Specifically, in the *visible-to-infrared* mode, NFS makes a marked improvement of 11.20% in Rank-1 and 8.64% in mAP compared to the top-performing method DDAG [66]. Similar increment also presents in the *infrared-to-visible* mode, which shows that our method is robust to multi-modal query settings. In conclusion, the above results clearly indicate the effectiveness of our automated feature search paradigm.

4.3. Ablation Study

This subsection studies the effectiveness of each module involved in NFS on SYSU-MM01 (*all* and *indoor search*, *single-shot* settings). As in Table 3, \mathcal{B} denotes the baseline two-stream network with the learning objective \mathcal{L}_b , \mathcal{N} represents the neural feature search block, and \mathcal{C} indicates the cross-modality contrastive optimization mechanism.

Table 3. Evaluation of each module on SYSU-MM01.

Method	All Search			Indoor Search		
	r1	r10	mAP	r1	r10	mAP
\mathcal{B}	47.00	84.11	46.46	52.70	89.30	60.93
$\mathcal{B} + \mathcal{N}$	48.91	86.03	47.92	54.12	92.20	62.31
$\mathcal{B} + \mathcal{C}$	52.29	90.22	50.91	57.23	94.14	65.32
$\mathcal{B} + \mathcal{N} + \mathcal{C}$	56.91	91.34	55.45	62.79	96.53	69.79

Effectiveness of Neural Feature Search. We evaluate how much improvement can be made by NFS with baseline learning objective \mathcal{L}_b . To be fair, all hyperparameters are fixed during evaluation. As shown in 2nd row of Table 3, NFS brings 1.91% Rank-1 and 1.46% mAP increases in *all search* mode compared with \mathcal{B} (row 2). Similar improvement can be observed in *indoor search* mode. This increment suggests that automated feature selection not only governs the baseline network to focus on informative parts of human bodies but also filters high sample noises.

Influence of Contrastive Optimization. Here, we investigate the contribution of contrastive loss. Considerable enhancement (5.29% of Rank-1 and 4.45% of mAP for *all search*, 4.53% of Rank-1 and 4.39% of mAP for *indoor search*) on the baseline model can be observed in Table 3.

This improvement manifests the superiority of contrastive loss for learning identity-related information. We further validate its effectiveness on NFS and the results are listed in the 4th row of Table 3. We observe that, with contrastive loss, NFS significantly surpasses the baseline model with 9.91% growth of Rank-1 and 8.99% gain of mAP for *all search*, while performance boost in *indoor search* is even more pronounced. Notably, from the comparison between 3rd and 4th row of Table 3, we see that NFS brings more benefits (4.62% of Rank-1 and 4.54% of mAP) to the baseline model with \mathcal{C} , demonstrating that the contrastive loss not only contributes to the optimization of \mathcal{B} , but also encourages NFS to discover more discriminative features.

Impact of Search Scope. In this part, we compare the performance of NFS conducted at different convolution stages (Table 4). The Rank-1 and mAP tend to increase when NFS is conducted on more layers. The best result appears when we perform NFS on Stage 1, 2, and 3. This is probably because, searching more stages expands the search space, which allows us to explore more varied selections of features. However, blindly extending the search scope will pose great challenges to the discovery of an optimal feature set. The underlying reason is that STE will generate more and more gradient estimation errors when backpropagating through too many layers [1].

Table 4. Comparison on NFS at different convolution stages.

	Stage1	Stage2	Stage3	Stage4	r1	mAP
1	✓	-	-	-	53.75	53.02
2	-	✓	-	-	53.41	52.67
3	-	-	✓	-	53.64	53.25
4	-	-	-	✓	53.92	52.83
5	✓	✓	-	-	55.62	54.31
6	✓	-	✓	-	54.85	53.97
7	✓	-	-	✓	54.45	53.76
8	-	✓	✓	-	55.13	54.07
9	-	✓	-	✓	54.41	53.29
10	-	-	✓	✓	54.42	53.87
11	✓	✓	✓	-	56.91	55.45
12	✓	✓	-	✓	55.21	53.13
13	✓	-	✓	✓	55.72	54.35
14	-	✓	✓	✓	55.94	54.62
15	✓	✓	✓	✓	55.18	54.21

4.4. Influence of Hyperparameters

In this subsection, we investigate the influence of hyperparameters involved in NFS, including the contrastive margin T (Eq. 12) and trade-off coefficient λ (Eq. 14). All results are based on SYSU-MM01 (*single-shot & all search*).

The Contrastive Margin T . Due to significant appearance differences between RGB and IR images, the original distance among negative pairs is relatively large. Thus, we tune the contrastive margin T from 10 to 20. The corresponding Rank-1 results are shown in Fig. 3(Left). NFS consistently outperforms the AGW baseline [67] with different margins and achieves best performance at $T = 15$.

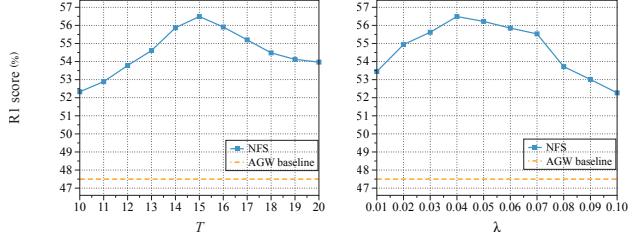


Figure 3. Parameter analysis for margin T and trade-off weight λ .

The Trade-off Coefficient λ . We also evaluate influence of the trade-off coefficient λ . Since the initial contrastive loss value may be very large, we consider λ from 0.01 to 0.1. As in Fig. 3(Right), consistent improvement can be observed again when we apply different λ . NFS achieves the best Rank-1 accuracy when $\lambda = 0.04$.

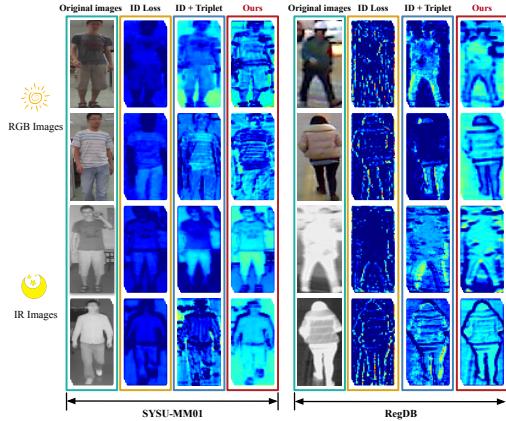


Figure 4. Visualization of feature maps produced by NFS with different loss functions on SYSU-MM01 and RegDB datasets. Best viewed in color.

4.5. Visualization of Learned Features

In order to inspect the effectiveness of our feature search based method, we visualize feature maps in the first shared block for 8 randomly selected images (4 samples per modality) on the two benchmark datasets (Fig. 4). It can be observed that, with the introduction of triplet loss (column 3), background noises are effectively eliminated while person information is preserved by search cells. We also see that significant improvement can be achieved when applying the contrastive loss (column 4) to NFS – not only irrelevant information is further filtered but also more discriminative cues are detected simultaneously.

We also examine the internal features captured by NFS using t-SNE [35]. As shown in Fig. 5, we visualize the learned representations of NFS and the baseline method on SYSU-MM01 and RegDB (5 randomly selected person identities per dataset). Specifically, Fig. 5(a) and 5(c) show the distribution of features extracted by the baseline method

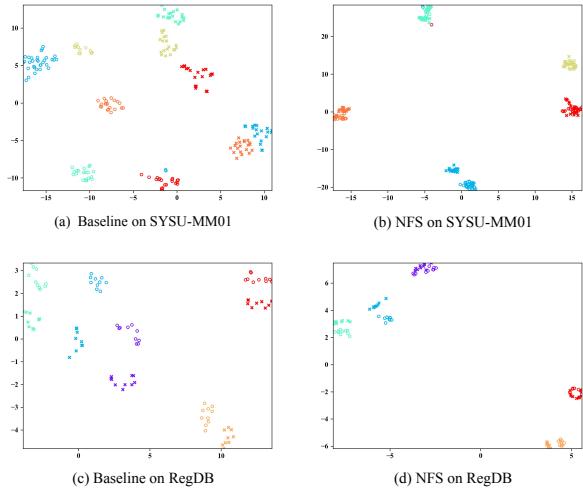


Figure 5. t-SNE visualization of the distribution of learned representations from NFS and the baseline method [67]. Different colors represent features of different identities. Circle and cross symbols refer to features of RGB and infrared images, respectively.

while Fig. 5(b) and 5(d) illustrate the NFS feature distribution. In comparison with Fig. 5(a) and 5(c), we see that feature distributions from visible and infrared modalities are fairly closer and less discriminable in Fig. 5(b) and 5(d). This indicates that NFS effectively minimizes the modality gap by aligning distributions of the two modalities. Furthermore, it is also observed that the proposed method separates feature points into disjoint clusters with larger inter-class margin while ensuring positive pairs from different modalities well aggregated. In a nutshell, NFS has a strong capability of detecting more discriminative cues in cross-modality settings.

5. Conclusion

This paper presents a novel insight of automated feature selection for RGB-IR ReID. A Neural Feature Search (NFS) paradigm is proposed to adaptively discover more identity characteristics. We first construct a dual-level feature search space, which makes it possible to jointly perform global-channel and local-spatial search operations. Then, we develop an efficient search algorithm to accelerate the selection process. Governed by a cross-modality contrastive optimization objective, this auto-searching algorithm is better able to select more high-quality invariant feature subsets for matching and retrieval. Experimental results on two standard RGB-IR ReID benchmarks demonstrate the effectiveness of NFS surpassing previous state-of-the-arts.

Acknowledgements. We are grateful to the AC panel and anonymous reviewers for their fruitful comments, corrections, and inspirations. We also thank Chuang Zhang, Yang Du, and Mengyao Tao for helpful discussions.

References

- [1] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013. 5, 7
- [2] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 5
- [3] Wuyang Chen, Xinyu Gong, Xianming Liu, Qian Zhang, Yuan Li, and Zhangyang Wang. Fasterseg: Searching for faster real-time semantic segmentation. *arXiv preprint arXiv:1912.10917*, 2019. 3
- [4] Seocheon Choi, Sumin Lee, Youngeun Kim, Taekyung Kim, and Changick Kim. Hi-cmd: Hierarchical cross-modality disentanglement for visible-infrared person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10257–10266, 2020. 2, 3, 4
- [5] Pingyang Dai, Rongrong Ji, Haibin Wang, Qiong Wu, and Yuyu Huang. Cross-modality person re-identification with generative adversarial training. In *IJCAI*, volume 1, page 2, 2018. 3, 6
- [6] Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. In *2005 IEEE computer society conference on computer vision and pattern recognition (CVPR'05)*, volume 1, pages 886–893. IEEE, 2005. 6
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. IEEE, 2009. 6
- [8] Changxing Ding, Kan Wang, Pengfei Wang, and Dacheng Tao. Multi-task learning with coarse priors for robust part-aware person re-identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 2
- [9] Thomas Elsken, Jan Hendrik Metzen, Frank Hutter, et al. Neural architecture search: A survey. *J. Mach. Learn. Res.*, 20(55):1–21, 2019. 3
- [10] Zhanxiang Feng, Jianhuang Lai, and Xiaohua Xie. Learning modality-specific representations for visible-infrared person re-identification. *IEEE Transactions on Image Processing*, 29:579–590, 2019. 3, 6
- [11] Mengxing Gong and Yijun Wang. An feature image generation based on adversarial generation network. In *2020 12th International Conference on Measuring Technology and Mechatronics Automation (ICMTMA)*, pages 479–482. IEEE, 2020. 1, 3
- [12] Yong Guo, Yaofu Chen, Yin Zheng, Peilin Zhao, Jian Chen, Junzhou Huang, and Mingkui Tan. Breaking the curse of space explosion: Towards efficient nas with curriculum search. In *International Conference on Machine Learning*, pages 3822–3831. PMLR, 2020. 3
- [13] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 1735–1742. IEEE, 2006. 5
- [14] Yi Hao, Nannan Wang, Jie Li, and Xinbo Gao. Hsme: Hypersphere manifold embedding for visible thermal person re-identification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8385–8392, 2019. 1, 2, 3, 6
- [15] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020. 5
- [16] 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. 6, 12
- [17] Xin He, Kaiyong Zhao, and Xiaowen Chu. Automl: A survey of the state-of-the-art. *Knowledge-Based Systems*, 212:106622, 2021. 2
- [18] Ruibing Hou, Bingpeng Ma, Hong Chang, Xinqian Gu, Shiguang Shan, and Xilin Chen. Interaction-and-aggregation network for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9317–9326, 2019. 1, 2
- [19] Mahdi M Kalayeh, Emrah Basaran, Muhittin Gökmén, Mustafa E Kamasak, and Mubarak Shah. Human semantic parsing for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1062–1071, 2018. 2
- [20] Jin Kyu Kang, Toan Minh Hoang, and Kang Ryoung Park. Person re-identification between visible and thermal camera images based on deep residual cnn using single input. *IEEE Access*, 7:57972–57984, 2019. 6
- [21] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in Neural Information Processing Systems*, 33, 2020. 5
- [22] Diangang Li, Xing Wei, Xiaopeng Hong, and Yihong Gong. Infrared-visible cross-modal person re-identification with an x modality. In *AAAI*, pages 4610–4617, 2020. 1, 3, 6
- [23] Zhihang Li, Teng Xi, Jiankang Deng, Gang Zhang, Shengzhao Wen, and Ran He. Gp-nas: Gaussian process based neural architecture search. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11933–11942, 2020. 3
- [24] Dongze Lian, Yin Zheng, Yintao Xu, Yanxiong Lu, Leyu Lin, Peilin Zhao, Junzhou Huang, and Shenghua Gao. Towards fast adaptation of neural architectures with meta learning. In *International Conference on Learning Representations*, 2019. 3
- [25] Feng Liang, Chen Lin, Ronghao Guo, Ming Sun, Wei Wu, Junjie Yan, and Wanli Ouyang. Computation reallocation for object detection. *arXiv preprint arXiv:1912.11234*, 2019. 3
- [26] Shengcui Liao, Yang Hu, Xiangyu Zhu, and Stan Z Li. Person re-identification by local maximal occurrence representation and metric learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2197–2206, 2015. 6

- [27] Hanxiao Liu, Karen Simonyan, Oriol Vinyals, Chrisantha Fernando, and Koray Kavukcuoglu. Hierarchical representations for efficient architecture search. In *International Conference on Learning Representations*, 2018. 3
- [28] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. In *International Conference on Learning Representations*, 2018. 2, 3, 4
- [29] Jinxian Liu, Bingbing Ni, Yichao Yan, Peng Zhou, Shuo Cheng, and Jianguo Hu. Pose transferable person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4099–4108, 2018. 1, 2
- [30] Jiawei Liu, Zheng-Jun Zha, Di Chen, Richang Hong, and Meng Wang. Adaptive transfer network for cross-domain person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7202–7211, 2019. 2
- [31] Gabriel Loaiza-Ganem and John P Cunningham. The continuous bernoulli: fixing a pervasive error in variational autoencoders. *arXiv preprint arXiv:1907.06845*, 2019. 2, 4
- [32] Yan Lu, Yue Wu, Bin Liu, Tianzhu Zhang, Baopu Li, Qi Chu, and Nenghai Yu. Cross-modality person re-identification with shared-specific feature transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13379–13389, 2020. 1, 2, 3, 6
- [33] Chuanchen Luo, Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. Spectral feature transformation for person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4976–4985, 2019. 2
- [34] 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. *IEEE Transactions on Multimedia*, 2019. 6
- [35] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008. 8
- [36] Jiaxu Miao, Yu Wu, Ping Liu, Yuhang Ding, and Yi Yang. Pose-guided feature alignment for occluded person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 542–551, 2019. 1, 2
- [37] Renato Negrinho and Geoff Gordon. Deeparchitect: Automatically designing and training deep architectures. *arXiv preprint arXiv:1704.08792*, 2017. 3
- [38] Dat Tien Nguyen, Hyung Gil Hong, Ki Wan Kim, and Kang Ryoung Park. Person recognition system based on a combination of body images from visible light and thermal cameras. *Sensors*, 17(3):605, 2017. 5
- [39] Juan-Manuel Pérez-Rúa, Valentin Vielzeuf, Stéphane Patteux, Moez Baccouche, and Frédéric Jurie. Mfas: Multimodal fusion architecture search. In *Proceedings of the IEEE Conference on computer vision and pattern recognition*, pages 6966–6975, 2019. 2, 3
- [40] Hieu Pham, Melody Guan, Barret Zoph, Quoc Le, and Jeff Dean. Efficient neural architecture search via parameters sharing. In *International Conference on Machine Learning*, pages 4095–4104. PMLR, 2018. 3, 4
- [41] Ruijie Quan, Xuanyi Dong, Yu Wu, Linchao Zhu, and Yi Yang. Auto-reid: Searching for a part-aware convnet for person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3750–3759, 2019. 2, 3
- [42] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V Le. Regularized evolution for image classifier architecture search. In *Proceedings of the aaai conference on artificial intelligence*, volume 33, pages 4780–4789, 2019. 3
- [43] Yantao Shen, Hongsheng Li, Shuai Yi, Dapeng Chen, and Xiaogang Wang. Person re-identification with deep similarity-guided graph neural network. In *Proceedings of the European conference on computer vision (ECCV)*, pages 486–504, 2018. 1, 2
- [44] Haotian Tang, Yiru Zhao, and Hongtao Lu. Unsupervised person re-identification with iterative self-supervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019. 1
- [45] Nihat Tekeli and Ahmet Burak Can. Distance based training for cross-modality person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019. 3
- [46] Maoqing Tian, Shuai Yi, Hongsheng Li, Shihua Li, Xuesen Zhang, Jianping Shi, Junjie Yan, and Xiaogang Wang. Eliminating background-bias for robust person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5794–5803, 2018. 1, 2
- [47] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*, 2019. 5
- [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 2
- [49] Guan'an Wang, Shuo Yang, Huanyu Liu, Zhicheng Wang, Yang Yang, Shuliang Wang, Gang Yu, Erjin Zhou, and Jian Sun. High-order information matters: Learning relation and topology for occluded person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6449–6458, 2020. 2
- [50] Guan'an Wang, Yang Yang, Jian Cheng, Jinqiao Wang, and Zengguang Hou. Color-sensitive person re-identification. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pages 933–939. AAAI Press, 2019. 2
- [51] Guan'an Wang, Tianzhu Zhang, Jian Cheng, Si Liu, Yang Yang, and Zengguang Hou. Rgb-infrared cross-modality person re-identification via joint pixel and feature alignment. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3623–3632, 2019. 1, 3, 5, 6, 7
- [52] Guan-An Wang, Tianzhu Zhang, Yang Yang, Jian Cheng, Jianlong Chang, Xu Liang, and Zeng-Guang Hou. Cross-modality paired-images generation for rgb-infrared person re-identification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 12144–12151, 2020. 1, 2, 3, 6
- [53] Linnan Wang, Rodrigo Fonseca, and Yuandong Tian. Learning search space partition for black-box optimization using monte carlo tree search. *arXiv preprint arXiv:2007.00708*, 2020. 3

- [54] Taiqing Wang, Shaogang Gong, Xiatian Zhu, and Shengjin Wang. Person re-identification by video ranking. In *European conference on computer vision*, pages 688–703. Springer, 2014. 1
- [55] Zhixiang Wang, Zheng Wang, Yinqiang Zheng, Yung-Yu Chuang, and Shin’ichi Satoh. Learning to reduce dual-level discrepancy for infrared-visible person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 618–626, 2019. 3, 6
- [56] Xing Wei, Diangang Li, Xiaopeng Hong, Wei Ke, and Yihong Gong. Co-attentive lifting for infrared-visible person re-identification. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 1028–1037, 2020. 1, 3
- [57] Ancong Wu, Wei-Shi Zheng, Shaogang Gong, and Jianhuang Lai. Rgb-ir person re-identification by cross-modality similarity preservation. *International journal of computer vision*, pages 1–21, 2020. 1
- [58] Ancong Wu, Wei-Shi Zheng, Hong-Xing Yu, Shaogang Gong, and Jianhuang Lai. Rgb-infrared cross-modality person re-identification. In *Proceedings of the IEEE international conference on computer vision*, pages 5380–5389, 2017. 1, 2, 3, 4, 5, 6
- [59] Bichen Wu, Xiaoliang Dai, Peizhao Zhang, Yanghan Wang, Fei Sun, Yiming Wu, Yuandong Tian, Peter Vajda, Yangqing Jia, and Kurt Keutzer. Fbnet: Hardware-aware efficient convnet design via differentiable neural architecture search. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10734–10742, 2019. 4
- [60] Guile Wu, Xiatian Zhu, and Shaogang Gong. Tracklet self-supervised learning for unsupervised person re-identification. In *AAAI*, pages 12362–12369, 2020. 1
- [61] Yichao Yan, Jie Qin, Bingbing Ni, Jiaxin Chen, Li Liu, Fan Zhu, Wei-Shi Zheng, Xiaokang Yang, and Ling Shao. Learning multi-attention context graph for group-based re-identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 2
- [62] Qize Yang, Ancong Wu, and Wei-Shi Zheng. Person re-identification by contour sketch under moderate clothing change. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. 2
- [63] Mang Ye, Xiangyuan Lan, Qingming Leng, and Jianbing Shen. Cross-modality person re-identification via modality-aware collaborative ensemble learning. *IEEE Transactions on Image Processing*, 29:9387–9399, 2020. 1, 2, 3, 6
- [64] Mang Ye, Xiangyuan Lan, Jiawei Li, and Pong C Yuen. Hierarchical discriminative learning for visible thermal person re-identification. In *AAAI*, 2018. 1, 3, 6
- [65] Mang Ye, Xiangyuan Lan, Zheng Wang, and Pong C Yuen. Bi-directional center-constrained top-ranking for visible thermal person re-identification. *IEEE Transactions on Information Forensics and Security*, 15:407–419, 2019. 2, 3, 6
- [66] Mang Ye, Jianbing Shen, David J Crandall, Ling Shao, and Jiebo Luo. Dynamic dual-attentive aggregation learning for visible-infrared person re-identification. *arXiv preprint arXiv:2007.09314*, 2020. 1, 2, 3, 4, 5, 6, 7
- [67] Mang Ye, Jianbing Shen, Gaojie Lin, Tao Xiang, Ling Shao, and Steven CH Hoi. Deep learning for person re-identification: A survey and outlook. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 1, 2, 4, 5, 6, 7, 8
- [68] Mang Ye, Jianbing Shen, and Ling Shao. Visible-infrared person re-identification via homogeneous augmented tri-modal learning. *IEEE Transactions on Information Forensics and Security*, 2020. 1, 2, 3
- [69] Mang Ye, Zheng Wang, Xiangyuan Lan, and Pong C Yuen. Visible thermal person re-identification via dual-constrained top-ranking. In *IJCAI*, volume 1, page 2, 2018. 1, 2, 3, 4, 6
- [70] Ziyue Zhang, Shuai Jiang, Congzhentao Huang, Yang Li, and Richard Yi Da Xu. Rgb-ir cross-modality person reid based on teacher-student gan model. *arXiv preprint arXiv:2007.07452*, 2020. 3
- [71] Liang Zheng, Yujia Huang, Huchuan Lu, and Yi Yang. Pose-invariant embedding for deep person re-identification. *IEEE Transactions on Image Processing*, 28(9):4500–4509, 2019. 2
- [72] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016. 2
- [73] Yang Zou, Xiaodong Yang, Zhiding Yu, BVK Kumar, and Jan Kautz. Joint disentangling and adaptation for cross-domain person re-identification. *arXiv preprint arXiv:2007.10315*, 2020. 1

The Architecture of Our Baseline Network

In our implementation, the baseline network adopts a commonly used two-stream structure (Fig. 5) and takes ResNet-50 as the backbone.

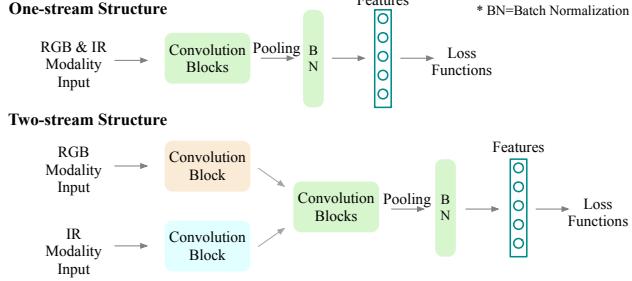


Figure 6. Two-stream structure of the baseline network.

Architecture details of the baseline network are shown in Table 5. All images are resized to 288×144 as the network inputs. The stride of the last convolutional block is set to 1 so as to obtain fine-grained feature maps. The other hyper-parameters are following [16] without tuning.

Table 5. Architecture details of our two-stream baseline network.

Layer name	Output size	50-layer	Type
conv1	144×72	$7 \times 7, 64$, stride 2	modality-specific
conv2	72×36	3×3 max pool, stride 2	shared
		$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	
conv3	36×18	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	shared
conv4	18×9	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	shared
conv5	18×9	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	shared

Visualization of Retrieved Examples

As shown in Fig. 7, the top-5 NFS retrieval results of 16 randomly selected query examples on the SYSU-MM01 dataset are plotted. We not only follow the original evaluation protocol, but also evaluate the *Visible-Infrared* setting.

In detail, the first column includes randomly selected samples from the query set, and retrieval results are sorted from left to right in descending order of cosine similarity scores. Due to lack of color information in IR images,

some cases are even difficult for human (e.g., query D and K). But the proposed method can retrieve correct results, which demonstrates the effectiveness of NFS in narrowing the large modality gap. According to the retrieval results for query B, C, D, and I, we also observe that our method exhibits certain robustness for high sample noises such as background clutter and partial occlusions. Interestingly, we discover that even if some persons change their clothes (e.g., query E), NFS can still return accurate retrieval results by mining other discriminative cues, perhaps the face or shoes. Another interesting phenomenon is that performance of *Visible-Infrared* setting is usually better than that of *Infrared-Visible* one. The main reason is that visible images often provide richer appearance information than their IR counterparts. Although there are still a few failure cases, most of these images (such as query F, O, and P) only present back views of the person-of-interest with limited identity-related information (e.g. face, texture, or logo of clothes). In conclusion, NFS exhibits promising performance in either query setting.



Figure 7. The top-5 retrieval results for 16 randomly selected query samples (8 samples per query setting) on the SYSU-MM01 dataset with our neural feature search method. Correct retrieved samples are in green boxes and wrong matchings are in red boxes (best viewed in color). Numerical values report cosine similarity scores of image pairs.