

Stronger Baseline for Person Re-Identification

Fengliang Qi, Bo Yan, Leilei Cao, Hongbin Wang
Ant Group

Abstract

Person re-identification (re-ID) aims to identify the same person of interest across non-overlapping capturing cameras, which plays an important role in visual surveillance applications and computer vision research areas. Fitting a robust appearance-based representation extractor with limited collected training data is crucial for person re-ID due to the high expense of annotating the identity of unlabeled data. In this work, we propose a **Stronger Baseline** for person re-ID, an enhancement version of the current prevailing method, namely, Strong Baseline, with tiny modifications but a faster convergence rate and higher recognition performance. With the aid of Stronger Baseline, we obtained the third place (i.e., 0.94 in mAP) in 2021 VIPriors Re-identification Challenge without the auxiliary of ImageNet-based pre-trained parameter initialization and any extra supplemental dataset.

1. Introduction

Person re-identification (re-ID) has become a core, and widely used technique in visual surveillance applications and computer vision research areas [23, 27, 17]. It aims at locating and recognizing a person of interest across multiple non-overlapping cameras in various spots [28, 22, 7].

The main challenge for the person-reID task is that the training dataset can be limited comparing with the conventional benchmark for classification task (e.g., ImageNet), making the risk of over-fitting increasing, which causes to the deterioration of the final generalization performance. To alleviate this dilemma, most of the current work focus on two main technical routines: 1) Designing customized and lightweight neutral network structure (e.g., OSNet); 2) Introducing discriminative margin-based loss function (e.g., Triplet loss) which has been widely used in metric learning and related research field (e.g., face verification).

In this work, we follow the second technical routine and propose the **Stronger Baseline** by modifying the Strong Baseline proposed by [19] with tiny modification but faster convergence rate and higher recognition performance as shown in Fig 1. In contrast to the argument claimed in [19],

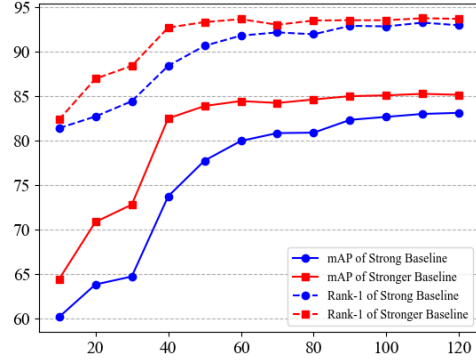


Figure 1. Convergence rate and performance are both enhanced considerably on Market1501 via the proposed Stronger Baseline.

we argue that the BNNeck is not the critical factor in alleviating the conflict between the triplet loss and cross-entropy loss minimization process. Instead, we claim the BNNeck is simply a standardization procedure. When combined with a softmax classifier, whose output follows the multinomial distribution, one of the exponential family distribution, the loss landscape can be smoother. The gradient update direction can be more stable, making the optimization algorithm more possible to arrive at the optimal global solution.

Furthermore, we attribute the occurrence of optimization conflict between cross-entropy loss and triplet loss to two main aspects: 1) inconsistency in identifying the hard samplers during the hard mining process and 2) inconsistency in gradient update direction during the objective minimizing process in two metric space (i.e., Cosine Metric Space and Euclidean Metric Space). As illustrated in Fig.2(a), compared to the anchor sample f_a , the positive sample f_p is easy in Cosine Metric Space since it lies in the similar radial direction with anchor sample f_a ; however, when it comes to the Euclidean Metric Space, the positive sample f_p should be treated as a hard sample since its significant Euclidean distance from anchor sample. The inconsistency in identifying the hardness of the positive sample holds for the negative sample (e.g., f_n in Fig.2(a)) in a similar principle as well. To alleviate this conflict, we use Batch-Normalization

(BN) module to standardize the distribution of the samplers as shown in Fig.2(b) and make the identification of samples' hardness consistent in both Cosine and Euclidean metric space.

As for the second inconsistency in gradient update direction, the gradient for the positive sample f_p in Euclidean Metric Space is parallel to the link between f_a and f_p . However, in Cosine Metric Space, the gradient update direction is parallel to the tangent direction, confusing the final comprehensive updating direction and deteriorating the optimization process. To mitigate this confusion, we use the **L2 Normalization operation** for the feature before triplet loss estimation. Thus the triplet loss can optimize the distance in the same cosine metric space with cross-entropy loss as illustrated in Fig.2(c).

Based on the analysis for alleviating the conflict in optimizing the cross-entropy loss and triplet loss, we propose our **Stronger Baseline** with tiny modification on the Strong Baseline proposed in [19]. As illustrated in Fig.3(b), the feature f_t extracted from the backbone network (e.g., ResNet) is processed by subsequent Batch-Normalization module to generate the feature f_i for cross-entropy loss calculation in the training stage and similarity evaluation in inference stage, which is the same as the pipeline of Strong baseline in Fig.3(a). Different from Strong Baseline, the **Stronger Baseline** calculate the triplet loss on the L2 normalization version of feature f_i instead of feature f_t . We discard the center loss for its limited influence on the final performance[19].

The main contribution of this work can be summarized in the following folds:

1. We observe two kinds of inconsistency when simultaneously optimizing the Cosine and Euclidean metric space, inducing the conflict in minimizing the cross-entropy and triplet loss.
2. We propose the **Stronger Baseline** based on analysis of how to alleviate the inconsistency during optimization by ameliorating the Strong Baseline with the limited modification.
3. We verify the superiority of our proposed method on the Market1501 and SynergySports benchmark.

2. Related Works

2.1. Person Re-identification

Gheissari et al. [8] first defined the person re-ID as a specific computer vision task. Most of the traditional methods focused on feature extraction via handcraft design before the emergence of deep learning[15, 6, 9, 20, 3]. With the rise of deep learning in recent years, convolutional neural networks (CNN) based feature representation methods have become the mainstream for image-based person re-ID [16, 27, 25, 26]. According to Zheng et al. [29], most of the

pre-existing image-based works concentrated on discriminative learning and metric learning. Discriminative learning [16, 1, 10, 18] aims at getting representative features for identity classification. While metric learning [24, 30] learns to project extracted features from different cameras and views into a common feature representation subspace.

2.2. Distance Metric Learning

The goals of the training and testing are slightly different in person re-ID tasks. The training process mainly focuses on classification or metric learning, while the testing process is a retrieval problem. Most existing supervised person re-ID methods apply identification loss for identity classification (e.g., cross-entropy loss)[24] and verification loss for metric learning (e.g., triplet loss)[12]. Triplet loss [21] is designed initially for face recognition problem and has been regarded as a commonly used method in the retrieval related tasks, especially in person re-ID. In the triplets, an anchor, a positive sample, and a negative sample are included. Since triplet loss is calculated by two randomly sampled person identification, it is difficult to ensure that the distance between the anchor and positive samples are smaller than the distance between the anchor and negative examples in the whole training dataset[19]. Quadruplet loss[2] is an improved version of triplet loss, which contains two different negative samples to learn a larger inter-class distance and a smaller intra-class distance compared with the triplet loss[13].

3. Our Method

3.1. Holistic Pipeline

In our experiments, the Stronger Baseline is applied as our holistic pipeline to realize the representation extracting and loss objective calculation. We use ResNet50 as a backbone for controlled experiments and *OSNet without ImageNet pretrain* for final result submission. The influence of the Backbone is empirically limited to the yielded performance, and the gain of the performance mostly attributes to the proposed **Stronger Baseline**. Both the BNNeck and FC layers are initialized through Kaiming initialization proposed in[11]. Each image in a training batch is preprocessed and resized to 256×128 in pixels. In each batch, we sample 8 identities, each with 4 images.

3.2. Data Augmentation

We utilize the following data augmentation for enhancing the generalization performance: *Random Horizontal Flip* with probability 0.5, *Random Erasing* with probability 0.5, *Color Jitter* with probability 0.5, and *AutoAugmentation*[4]. We are inconclusive in which augmentation strategy is beneficial to the final performance at most. We solely

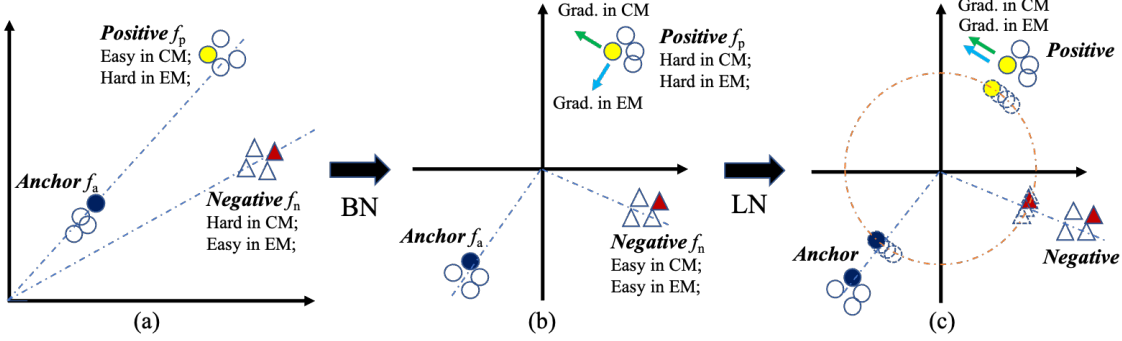


Figure 2. **Schematic illustration of the optimization conflict in Cosine and Euclidean metric space.** (a) The hardness of samples is identified inconsistently in Cosine and Euclidean metric space. (b) BN module can alleviate the first inconsistency in identifying the hardness of samples. (c) LN operation can alleviate the second inconsistency in gradient updating direction. Shape refers to the identity of the samples, and filling color refers to the role in metric pairs (e.g., anchor), Grad., BN, LN, CM, and EM refer to the Gradient, Batch-Normalization, L2 Normalization, Cosine Metric, and Euclidean Metric, respectively.

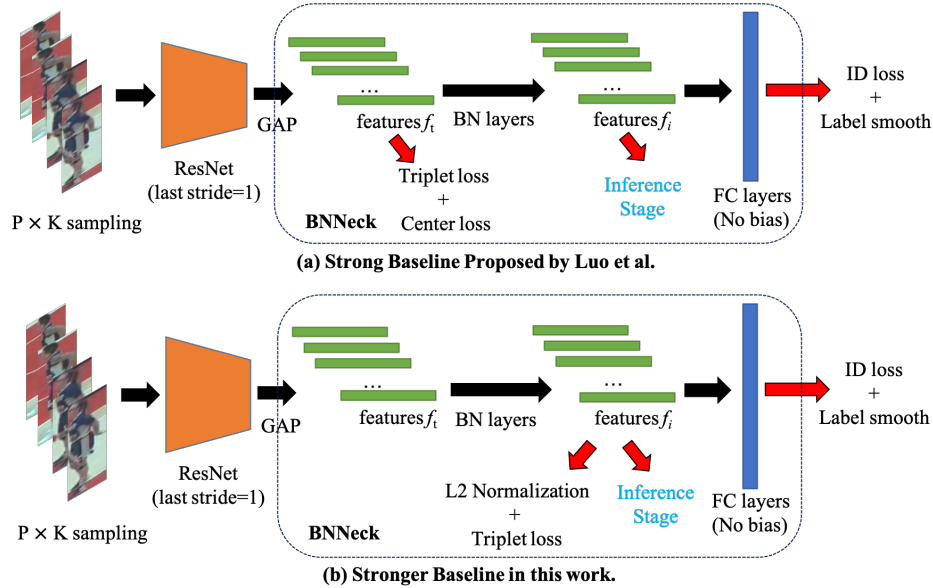


Figure 3. **The overview of the pipeline of Stronger Baseline.** (a) The primary pipeline of Strong Baseline. (b) The enhanced version of Strong Baseline, namely, Stronger Baseline.

adopted the first two augmentation strategies while conducting the controlled experiment.

3.3. Post-processing

Person re-ID can also be regarded as a retrieval problem. While presenting our final results, we utilize two common post-processing strategies, i.e., **Query Expansion** and **ReRank**, whose hyperparameters follow the default setting in FastReID.

3.4. Optimization Strategy

For controlled trials on Market1501, Adam[14] optimizer is selected with an initial learning rate of 3.5×10^{-4} . The commonly adopted warm-up strategy[5] is applied to

bootstrap the network for better performance. In practice, the network is optimized for 120 epochs. We spend 10 epochs linearly increasing the learning rate from 3.5×10^{-6} to 3.5×10^{-4} , and it then decays at the 30th and 55th epoch. We discard the warm-up stage for final submission and use SGD optimizer instead due to the exclusion of pretrain. The network is optimized for 350 epochs, and the init learning rate of 0.065 decays at the 150th, 225th and 300th epoch by 0.1.

4. Experiments

We conduct two main experiments in this work, i.e., Controlled Experiment on the public benchmark (i.e., Market1501) and Submitted Experiment on SynergySports. The

former aims to justify the superiority of the proposed Stronger Baseline, and the latter aims to obtain a better result with extra tricks (e.g., Complex Data Augmentation and Post-Processing Strategy).

4.1. Datasets and Evaluation Protocol

Market-1501 dataset is featured by 1,501 IDs, 19,732 gallery images and 12,936 training images captured by 6 cameras. Market-1501 are produced by the DPM detector. The Cumulative Matching Characteristics (CMC) curve is used for performance evaluation, which encodes the possibility that the query person is found within the top n ranks in the rank list. We also employ the mean Average Precision (mAP), which considers the retrieval process's precision and recall. The evaluation toolbox provided by the Market-1501 authors is used.

SynergySports generates from short sequences of basketball games, and each sequence is composed of 20 frames. For the validation and test sets, the query images are persons taken at the first frame, while the gallery images are identities taken from the 2nd to the last frame. There are 436 IDs, 8569 images in training split, and 50 IDs, 960 images in validation split.

4.2. Empirical Results

As illustrated in Fig 1, the Stronger Baseline can enhance the convergence rate and performance simultaneously. The detailed results are listed in Table 1. The mAP of Strong Baseline is improved by 2.03, which is a considerable enhancement considering the tiny modification and no extra introduced overhead. We further evaluate our method on the validation set of SynergySports; the results are listed in Table 2, the mAP achieves the predictable improvement by 1.89.

As for the Submitted Experiment, we switch to the OS-Net, a lightweight backbone customized for person re-ID, and introduce the complex data augmentation and post-processing strategy for better generalization performance. The results on the validation and test set are listed in Table 3.

Protocol	Strong Baseline	Stronger Baseline
mAP	83.13	85.16 (+2.03)
Rank-1	92.99	93.71 (+0.72)

Table 1. The mAP and rank-1 comparison result of Strong Baseline and Stronger Baseline on Market1501. We use ResNet50 as backbone and ignore the complex data augmentation (i.e., Color Jitter and AutoAug) and Post-Processing here.

5. Conclusions

In this paper, we concluded two kinds of inconsistency while optimizing the Cosine and Euclidean metric space

Protocol	Strong Baseline	Stronger Baseline
mAP	92.74	94.63 (+1.89)

Table 2. The mAP comparison result of Strong Baseline and Stronger Baseline on SynergySports (Validation Set). We use ResNet50 as backbone and ignore the complex data augmentation (i.e., Color Jitter and AutoAug) and Post-Processing here.

Protocol	Validation Set	Test Set
mAP	95.17	94.19

Table 3. The mAP result on SynergySports. We use OSNet1x0 as backbone and introduce the complex data augmentation (i.e., Color Jitter and AutoAug) and Post-Processing here.

simultaneously, i.e., hardness identification inconsistency and gradient update inconsistency, causing the conflict between minimizing the cross-entropy loss and triplet loss. To alleviate the inconsistency and ameliorate the optimization process, we proposed the Stronger Baseline with tiny modifications on the Strong Baseline but a faster convergence rate and higher evaluation performance. With the aid of Stronger Baseline, we obtain a third place in the 2021 VIPriors Re-identification Challenge without the auxiliary of ImageNet-based pre-train parameter initialization and any extra supplemental dataset.

References

- [1] Xiaobin Chang, Timothy M Hospedales, and Tao Xiang. Multi-level factorisation net for person re-identification. In *CVPR*, pages 2109–2118, 2018. 2
- [2] Weihua Chen, Xiaotang Chen, Jianguo Zhang, and Kaiqi Huang. Beyond triplet loss: A deep quadruplet network for person re-identification. In *CVPR*, 2017. 2
- [3] Dong Seon Cheng, Marco Cristani, Michele Stoppa, Loris Bazzani, and Vittorio Murino. Custom pictorial structures for re-identification. In *BMVC*, volume 1, page 6, 2011. 2
- [4] Ekin Dogus Cubuk, Barret Zoph, Dandelion Mané, Vijay Vasudevan, and Quoc V. Le. Autoaugment: Learning augmentation policies from data. *CoRR*, abs/1805.09501, 2018. 2
- [5] Xing Fan, Wei Jiang, Hao Luo, and Mengjuan Fei. Sphered: Deep hypersphere manifold embedding for person re-identification. *Journal of Visual Communication and Image Representation*, 60:51–58, 2019. 3
- [6] Michela Farenzena, Loris Bazzani, Alessandro Perina, Vittorio Murino, and Marco Cristani. Person re-identification by symmetry-driven accumulation of local features. In *CVPR*, pages 2360–2367, 2010. 2
- [7] Yang Fu, Xiaoyang Wang, Yunchao Wei, and Thomas Huang. Sta: Spatial-temporal attention for large-scale video-based person re-identification. In *AAAI*, volume 33, pages 8287–8294, 2019. 1
- [8] Niloofar Gheissari, Thomas B Sebastian, and Richard Hartley. Person reidentification using spatiotemporal appearance. In *CVPR*, pages 1528–1535, 2006. 2
- [9] Douglas Gray and Hai Tao. Viewpoint invariant pedestrian recognition with an ensemble of localized features. In *ECCV*, pages 262–275, 2008. 2

- [10] Yiluan Guo and Ngai-Man Cheung. Efficient and deep person re-identification using multi-level similarity. In *CVPR*, pages 2335–2344, 2018. 2
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *ICCV*, pages 1026–1034, 2015. 2
- [12] Alexander Hermans, Lucas Beyer, and Bastian Leibe. In defense of the triplet loss for person re-identification. *arXiv:1703.07737*, 2017. 2
- [13] Ming Jiang, Biao Leng, Guanglu Song, and Zhijun Meng. Weighted triple-sequence loss for video-based person re-identification. *Neurocomputing*, 381:314–321, 2020. 2
- [14] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014. 3
- [15] Igor Kviatkovsky, Amit Adam, and Ehud Rivlin. Color invariants for person reidentification. *IEEE Transactions on pattern analysis and machine intelligence*, 35(7):1622–1634, 2012. 2
- [16] Wei Li, Xiatian Zhu, and Shaogang Gong. Harmonious attention network for person re-identification. In *CVPR*, pages 2285–2294, 2018. 2
- [17] Chih-Ting Liu, Chih-Wei Wu, Yu-Chiang Frank Wang, and Shao-Yi Chien. Spatially and temporally efficient non-local attention network for video-based person re-identification. In *arXiv:1908.01683*, 2019. 1
- [18] Hao Liu, Jiashi Feng, Meibin Qi, Jianguo Jiang, and Shuicheng Yan. End-to-end comparative attention networks for person re-identification. *IEEE Transactions on Image Processing*, 26(7):3492–3506, 2017. 2
- [19] Hao Luo, Youzhi Gu, Xingyu Liao, Shenqi Lai, and Wei Jiang. Bag of tricks and a strong baseline for deep person re-identification. In *CVPRW*, 2019. 1, 2
- [20] Tetsu Matsukawa, Takahiro Okabe, Einoshin Suzuki, and Yoichi Sato. Hierarchical gaussian descriptor for person re-identification. In *CVPR*, pages 1363–1372, 2016. 2
- [21] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, pages 815–823, 2015. 2
- [22] Arulkumar Subramaniam, Athira Nambiar, and Anurag Mitral. Co-segmentation inspired attention networks for video-based person re-identification. In *ICCV*, pages 562–572, 2019. 1
- [23] Han Sun, Zhiyuan Chen, Shiyang Yan, and Lin Xu. Mvp matching: A maximum-value perfect matching for mining hard samples, with application to person re-identification. In *ICCV*, pages 6737–6747, 2019. 1
- [24] Yifan Sun, Liang Zheng, Weijian Deng, and Shengjin Wang. Svdnet for pedestrian retrieval. In *ICCV*, pages 3800–3808, 2017. 2
- [25] 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. 2
- [26] Tong Xiao, Hongsheng Li, Wanli Ouyang, and Xiaogang Wang. Learning deep feature representations with domain guided dropout for person re-identification. In *CVPR*, pages 1249–1258, 2016. 2
- [27] Shiyang Yan, Jun Xu, Yuai Liu, and Lin Xu. Hornet: a hierarchical offshoot recurrent network for improving person re-id via image captioning. *arXiv:1908.04915*, 2019. 1, 2
- [28] Ruimao Zhang, Jingyu Li, Hongbin Sun, Yuying Ge, Ping Luo, Xiaogang Wang, and Liang Lin. Scan: Self-and-collaborative attention network for video person re-identification. *IEEE Transactions on Image Processing*, 28(10):4870–4882, 2019. 1
- [29] Liang Zheng, Yi Yang, and Alexander G Hauptmann. Person re-identification: Past, present and future. *arXiv:1610.02984*, 2016. 2
- [30] Liang Zheng, Hengheng Zhang, Shaoyan Sun, Manmohan Chandraker, Yi Yang, and Qi Tian. Person re-identification in the wild. In *CVPR*, pages 1367–1376, 2017. 2