

# Compact Network Training for Person ReID

Hussam Lawen\*

Avi Ben-Cohen\*

hussam.lawen@alibaba-inc.com

avi.bencohen@alibaba-inc.com

DAMO Academy, Alibaba Group

Tel-Aviv, Israel

Itamar Friedman

itamar.friedman@alibaba-inc.com

DAMO Academy, Alibaba Group

Tel-Aviv, Israel

Matan Protter

matan.protter@alibaba-inc.com

DAMO Academy, Alibaba Group

Tel-Aviv, Israel

Lihl Zelnik-Manor

lihi.zelnik@alibaba-inc.com

DAMO Academy, Alibaba Group

Tel-Aviv, Israel

## ABSTRACT

The task of person re-identification (ReID) has attracted growing attention in recent years leading to improved performance, albeit with little focus on real-world applications. Most SotA methods are based on heavy pre-trained models, e.g. ResNet50 (~25M parameters), which makes them less practical and more tedious to explore architecture modifications. In this study, we focus on a small-sized randomly initialized model that enables us to easily introduce architecture and training modifications suitable for person ReID. The outcomes of our study are a compact network and a fitting training regime. We show the robustness of the network by outperforming the SotA on both Market1501 and DukeMTMC. Furthermore, we show the representation power of our ReID network via SotA results on a different task of multi-object tracking.

## KEYWORDS

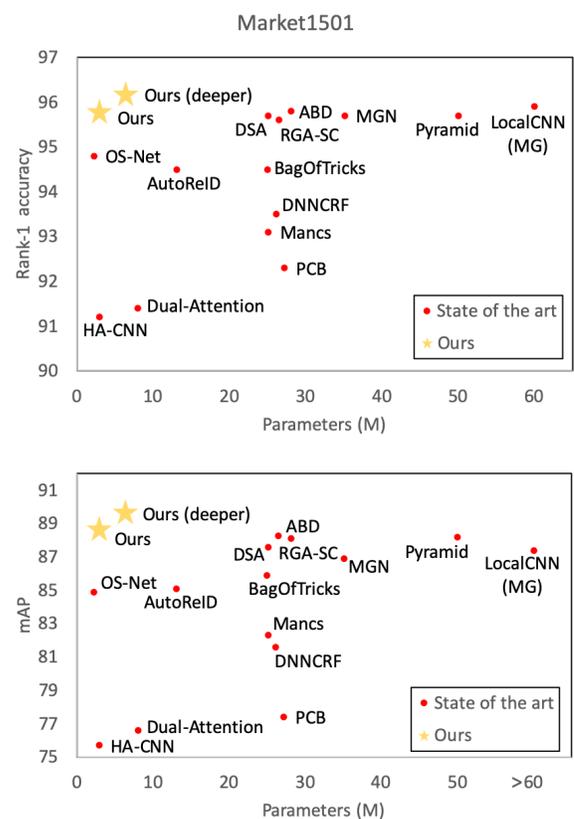
Deep person ReID, multi-object tracking, compact network

## 1 INTRODUCTION

The objective in person re-identification (ReID) is to assign a stable ID to a person in multiple camera views. In this study we are interested in the development of small sized models for ReID with high accuracy for two main reasons. First, it is beneficial for practical deployment and productization of ReID solutions. Second, the research for models that provide high accuracy requires exploration of many architecture variations and training schemes. When the backbone is heavy, re-training consumes both a lot of time and computing resources which we wish to avoid. Our approach differs from many state-of-the-art (SotA) methods, that rely on large pre-trained backbone models, such as ResNet50, e.g. [11, 23, 25, 27].

We argue that a cost-effective ReID model should be computationally efficient, capable of running on low-res video input, and robust to multiple camera setting. Hence, we propose an efficient ReID model and training schemes that demonstrate state of the art performance under these requirements. To reduce the computational burden, we aim to decrease the number of parameters and use a relatively small ReID model. Figure 1 shows the current state of the art results [1, 2, 10, 11, 16, 21, 23, 25, 26, 28, 30–32, 38] and the number of parameters compared to our proposed method on

\*Both authors contributed equally to this research.



**Figure 1: Performance comparison of our approach and SotA ReID methods on Market1501 dataset. Top: rank-1 accuracy vs. number of parameters. Bottom: mAP vs. number of parameters.**

the popular Market1501 dataset [33] in terms of rank-1 accuracy and mAP. For some methods, the number of parameters was not known so we used an estimated lower bound. Using our proposed training framework we achieve state of the art results with an order of magnitude smaller model compared to the best existing ReID CNN.

The importance of training “tricks” for deep person ReID has been discussed before in [11]. In this paper, we suggest training techniques and architecture modifications that improve the harmonious attention network HA-CNN of [10] to achieve similar or better results than much larger and complicated models. The contribution of this paper is thus three-fold:

- We propose a compact and robust deep person ReID model. Our model achieves state of the art results on two popular person ReID datasets (Market1501 and DukeMTMC ReID [18]). This is despite having a small number of parameters, small number of FLOPS, and low resolution input image, in comparison to other leading methods.
- We study a variety of training schemes and network choices that prove useful. While we have explored their affect only for HA-CNN, we believe they could be of interest for others to examine in other setups.
- We demonstrate the utility of the proposed person ReID model also for other tasks, by improving multi-target multi-camera tracking.

In the following section we describe the baseline ReID network we started with. The training techniques and architecture modifications that were explored in this study are presented in section 3. Next, the experimental results including an ablation study, additional analysis, and comparison to state of the art are presented (section 4). Finally, multi camera multi target tracking results are presented in section 5.

## 2 BASELINE REID NETWORK - HA-CNN

Our goal is a compact model that gives high accuracy with low-resolution input images, in order to reduce computational complexity. Therefore, we chose as baseline the light-weight Harmonious Attention CNN (HA-CNN) [10]. HA-CNN is sufficiently compact to be trained from scratch thus obviating the need to pre-train on additional data. Nonetheless, it provides good results taking into consideration its small number of parameters (2.7M). In addition, the input image size for this network is relatively small compared to other person ReID networks.

HA-CNN is an attention network with several attention modules including soft spatial and channel-wise attention and hard attention to extract local regions. The network architecture holds two branches: a global one, and a local one that uses the regions extracted based on the hard attention. Finally, the output vectors of both branches are concatenated for the final person image descriptor. Holding two branches and multiple attention modules improves the network perception, despite these features the HA-CNN keeps a small number of parameters making it accurate and efficient. However, parts of the architecture can still be optimized as well as the training scheme. Optimizing it can further improve the HA-CNN and obtain a more accurate model.

## 3 METHODS

It is well known that the performance of deep learning models is highly dependent on both the choice of architecture and the training scheme. Specifically, recent work has shown that training procedure refinements can significantly improve ReID results [11]. In the following we explore training schemes (Section 3.1) and architecture

modifications (section 3.2) that lead to better ReID performance of HA-CNN. To make our survey more complete we further mention several modifications that did not improve the model performance (Section 3.3).

### 3.1 Training techniques

The following training techniques were found useful by our study:

*Weighted triplet loss with Soft margin.* The triplet loss is widely used to train Person ReID models, as well as other computer vision tasks such as Face Recognition and Few-Shot Learning. The original triplet loss was proposed by Schroff *et al.* [20]. We denote an anchor sample by  $x_a$ , positive samples as  $x_p \in P(a)$  and negative samples as  $x_n \in N(a)$ , then the triplet loss can be written as:

$$L_1 = [m + d(x_a, x_p) - d(x_a, x_n)]_+ \quad (1)$$

where  $m$  is the given inter-class separation margin,  $d$  denotes distance of appearance, and  $[\cdot]_+ = \max(0, \cdot)$ .

Hermans *et al.* [6] proposed the batch-hard triplet loss that selects only the most difficult positive and negative samples:

$$L_2 = \left[ m + \max_{x_p \in P(a)} d(x_a, x_p) - \min_{x_n \in N(a)} d(x_a, x_n) \right]_+ \quad (2)$$

In contrast to the original triplet loss, the batch-hard triplet loss emphasizes hard examples. However, it is sensitive to outlier samples and may discard useful information due to its hard selective approach. To deal with these problems, Ristani *et al.* proposed the batch-soft triplet loss:

$$L_3 = \left[ m + \sum_{x_p \in P(a)} w_p d(x_a, x_p) - \sum_{x_n \in N(a)} w_n d(x_a, x_n) \right]_+ \quad (3)$$

$$w_p = \frac{e^{d(x_a, x_p)}}{\sum_{x \in P(a)} e^{d(x_a, x)}}, \quad w_n = \frac{e^{-d(x_a, x_n)}}{\sum_{x \in N(a)} e^{-d(x_a, x)}}$$

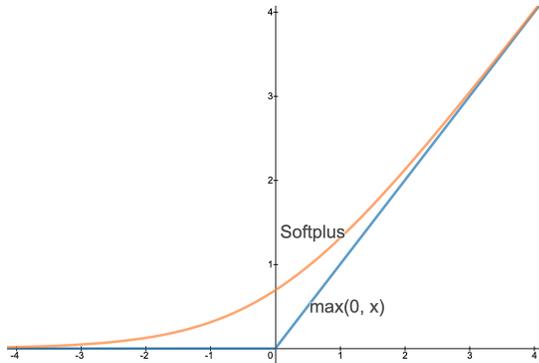
Observe, that the hyper-parameter  $m$ , which denotes the margin, exists in all of these triplet loss variations. Tuning this hyper-parameter manually is not easy, therefore, we next propose an alternative triplet loss that eliminates it.

Our key idea is to replace the hard cutoff max function with an exponential decay *Softplus*( $\cdot$ ) =  $\ln(1 + \exp(\cdot))$  as follows:

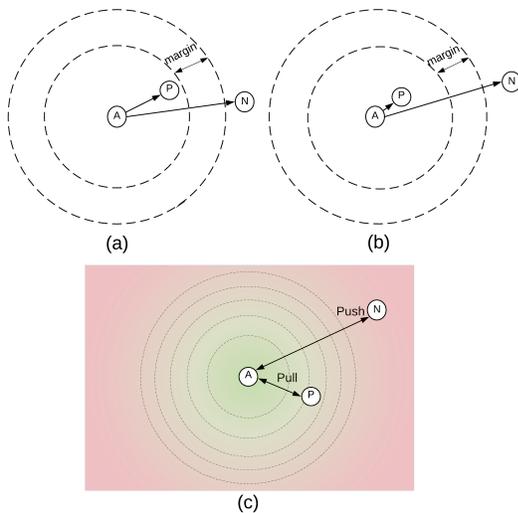
$$L_4 = \text{Softplus} \left( \sum_{x_p} w_p d(x_a, x_p) - \sum_{x_n} w_n d(x_a, x_n) \right) \quad (4)$$

The soft margin eliminates the margin parameter.

Figure 3 illustrates one of the benefits of the soft margin over the hard margin. Using a hard margin value, when the separation between the negative samples and the positive samples becomes larger than the hard margin, the loss is zero and therefore further minimization will not push the positive samples closer or the negative samples farther away from the anchor. This is illustrated in the examples in (a) and (b) that will both obtain a loss value of zero since both answer the assumption of the hard margin. Conversely, the soft margin encourages a continuous reduction of the positive distance to the anchor while increasing the negative distance. This is illustrated in (c), that shows the the computed loss will continue



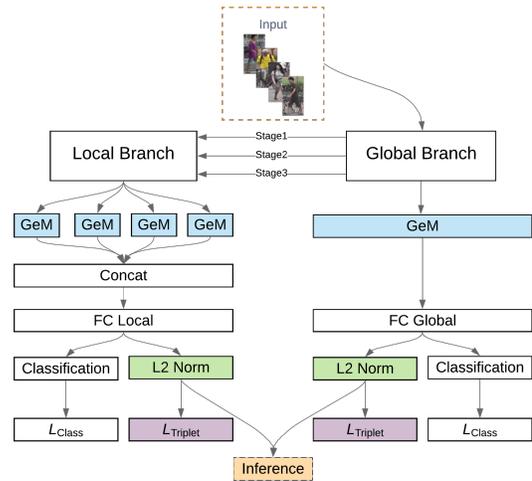
**Figure 2: The Softplus function ( $\ln(1 + \exp(\cdot))$ ) compared to  $\max(0, \cdot)$ .**



**Figure 3: Example of hard margin vs soft margin. The scenario in (b) is more desirable than (a) because the positive sample is closer to the anchor and the negative sample is farther away. This, however, will not be captured by the hard margin triplet loss because both cases correspond to a loss value of zero. (c) Differently, when using a soft margin the loss will continue to pull the positive sample closer to the anchor while pushing the negative sample away and will encourage going from (a) to (b).**

to push the positive sample closer to the anchor while pushing the negative sample away.

*L2 normalization.* The normalization of the feature vectors can be important when using two different loss functions such as cross-entropy and triplet loss which are optimized using different distance measures. [11] tackled the normalization problem by adding a batch normalization layer after the feature vectors, right before the fully connected layer. In our empirical studies we found that simply using  $L_2$  normalization for each feature vector (global and local) during training achieves an even better performance. Figure 4 shows the additional  $L_2$  normalization used during training and inference.



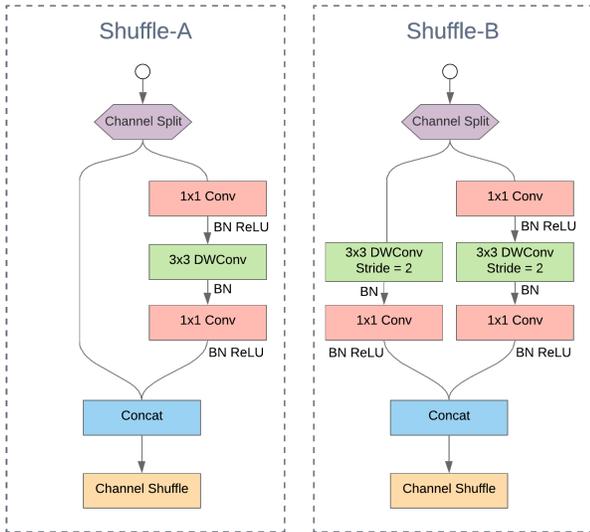
**Figure 4: Our ReID architecture shows the proposed modifications over the original HA-CNN:  $L_2$  normalization during training, GeM instead of average pooling, and soft triplet loss.**

*SWAG [13].* A common technique to further boost the performance of a model is via ensembles. A common approach is to use an ensemble of models in test time for the final prediction, however, this requires high computing resources. A more efficient approach is Stochastic weight averaging (SWA) [8], that forms an ensemble during training and outputs a single model for inference. SWA essentially conducts a uniform average over several model weights traversed by SGD during training to achieve a wider region of the loss minima. In order to use SWA a learning rate scheduler is required.

We have made two modifications over HA-CNN ensemble scheme. First, we follow [13] and use SWA-Gaussian (SWAG). SWAG fits a Gaussian distribution using the SWA solution and diagonal covariance forming an approximate posterior distribution over neural network weights. Next, SWAG performs a Bayesian model averaging based on the Gaussian distribution. Second, we have found empirically that the original learning rate scheduler of [8] can be improved. We suggest using the cosine annealing learning rate scheduler with  $cycles = 15$  of 35 epochs and cycle decay factor of 0.7 after each cycle. At the end of each cycle we average the weights of the current model with the previous models taken from the end of each cycle.

*Other training techniques.* The random erasing augmentation (REA) [36] that randomly erases a rectangle in an image has shown to improve the model generalization ability. We used REA with the following parameters: probability for random erasing an image of 0.5, area ratio of erasing a region in the range of  $0.02 < S_e < 0.4$ , and with aspect ratio in the range of  $0.3 < r < 3.3$ .

*Warmup [4]* - used to bootstrap the network for better performance. Starting with a smaller learning rate has shown to improve the training process stability, especially when using a randomly initialized model. Using warmup we start the training with a small learning rate and then gradually increase it. We used the following



**Figure 5: The shuffle blocks used in this study to replace the original HA-CNN inception blocks.**

learning rate scheme:

$$lr(t) = \begin{cases} 3 \times 10^{-2} \times \frac{t}{10} & \text{if } t \leq 10 \\ 3 \times 10^{-2} & \text{if } 10 < t \leq 150 \\ 3 \times 10^{-3} & \text{if } 150 < t \leq 225 \\ 3 \times 10^{-4} & \text{if } 225 < t \leq 350 \end{cases} \quad (5)$$

Label smoothing [24] - widely used for classification problems by encouraging the model to be less confident during training and prevent over-fitting. We used label smoothing in a similar way as proposed in [11].

### 3.2 Architecture modifications

In addition to the training techniques listed above, we further suggest the following architecture modifications to HA-CNN.

*Shuffle blocks [12].* Our goal was to improve the network accuracy while maintaining a small number of parameters. To do this, we examined replacing the inception blocks with the shuffle blocks presented in Figure 5.

Shuffle-A is more efficient than the original inception block since it splits the input features into two equal branches, the first branch remains as is while three convolution operators are applied to the second branch. In addition, one of the convolution operators is depth-wise convolution. The Shuffle-A block can be used in a repeated sequence and still maintain the same number of parameters as the original inception block. Hence, we were able to build a deeper network with a similar number of parameters. The Shuffle-B block is similar to Shuffle-A but can be used for spatial down-sampling or channel expansion. These characteristics require convolution operators to be applied also to the first branch. Table 1 summarizes the repeated sequences of Shuffle blocks used in our proposed architecture.

*Generalized Mean (GeM) [17].* In the original HA-CNN global average pooling was used just before the fully connected layer. Replacing it with global max pooling gave undecive results, sometime better and sometimes worse. Therefore, we suggest using the trainable Generalized Mean (GeM) pooling, which generalizes both max and average pooling. The GeM operator for a single feature map  $f_k$  can be written as:

$$GeM(f_k = [x_0, x_1, \dots, x_n]) = \left[ \frac{1}{n} \sum_{i=1}^n x_i^{p_k} \right]^{\frac{1}{p_k}} \quad (6)$$

We initialized the parameter  $p_k = 3$ . Figure 4 shows where it is used during training and inference.

*Deeper and wider.* We further study empirically the impact of using a deeper and wider version of the architecture by modifying the number of shuffle blocks as well as the number of output channels in each stage. Table 1 presents these modifications in bold.

### 3.3 Additional tricks we tried

For completeness, we list here training options that have been introduced by prior work and our experiments found to deteriorate the results:

- (1) As mentioned before, max and average pooling provide different results so one way to benefit from both pooling methods is by concatenation of their output. Basically we tried to replace the global average pooling used in the original HA-CNN architecture with these two pooling methods and concatenations. It resulted in a similar accuracy with more parameters in the final model.
- (2) The batch norm suggested by [11] provided inferior results when compared to the simple  $L_2$  normalization.
- (3) Hard triplet loss instead of the soft version was too sensitive to outliers.
- (4) Shuffle blocks without  $L_2$  normalization or soft margin in the triplet loss didn't improve the performance.
- (5) Training for more epochs didn't improve the performance. The only way it did lead to an improvement was using the Cyclic LR scheme.
- (6) Cyclic LR scheme didn't improve the results when used from scratch from the beginning of the training. It only worked when used in additional training epochs after the the model converged.

## 4 EXPERIMENTAL RESULTS

In the following we evaluate our models on Market1501 and DukeMTMC ReID datasets based on rank-1 accuracy and mAP. Next, the performance boost by each methods presented in section 3 is evaluated.

*Implementation details.* All person images are resized to  $160 \times 64$ . We used SGD for optimization with a linear warm-up as in Equation (5) for a total of 350 epochs. When using SWAG we train for 15 cycles of 35 epochs which sums up to 525 additional epochs. We randomly sample 8 identities and 4 images per person in each training batch.

Local Branch	Global Branch	Layer	Input	Stride	1×		2×	
					Repeat	Output Ch.	Repeat	Output Ch.
	Conv1	Conv 3x3	160×64	2	1	32	1	<b>36</b>
	Stage1	Shuffle-B	80×32	1	1	128	1	<b>240</b>
		Shuffle-A		1	7		<b>8</b>	
		Shuffle-B		2	1		1	
	Soft-Attn1	HA-Block	40×16	1	1		1	
Hard-Attn1		Shuffle-B	4×(24×28)	1	1		1	
	Stage2	Shuffle-B	40×16	1	1	256	1	<b>320</b>
		Shuffle-A		1	10		<b>11</b>	
		Shuffle-B		2	1		1	
	Soft-Attn2	HA-Block	20×8	1	1		1	
Hard-Attn2		Shuffle-B	4×(12×14)	1	1		1	
	Stage3	Shuffle-B	20×8	1	1	384	1	<b>480</b>
		Shuffle-A		1	7		<b>8</b>	
		Shuffle-B		2	1		1	
	Soft-Attn3	HA-Block	10×4	1	1		1	
Hard-Attn3		Shuffle-B	4×(6×7)	1	1		1	
	Pooling	GeM	10×4	1	1		1	
Pooling		GeM	4×(3×4)	1	1		1	
	FC Global	Linear	1×1	1	1	512	1	<b>960</b>
FC Local		Linear	1×1	1	1		1	
FLOPs					0.72B		<b>1.68B</b>	
# of Params.					2.9M		<b>6.4M</b>	

**Table 1: Overall architecture of our model, for 2 different levels of complexities. Since our architecture uses a low-resolution input of 160x64, we down-scale the feature maps by applying strided convolution only in the last layer of each stage and not in the beginning. This way the network can leverage a higher spatial resolution in most of the network.**

#### 4.1 Comparison to state of the arts

We compare our models performance to several state of the art methods (Table 2). Our best model achieves state of the art results in terms of rank-1 accuracy and mAP on Market1501 (96.2, 89.7) and DukeMTMC (89.8, 80.3) with only 6.4M parameters. To our best knowledge, our model achieves the best performance on these public datasets. It should be noted that the smaller version of our model (2.9M parameters) also achieves state of the art results on both datasets.

In terms of FLOPS our final network has 1.7B FLOPS while the ResNet 50 used in Luo *et al.* [11] implementation has 4.1B FLOPS. We did not apply re-ranking for clear comparison and since it is currently not relevant for real world practice.

#### 4.2 Ablation study

To evaluate the different training techniques explored in this study we set several experiments in an ablation study. Table 3 shows the different modifications starting from the original HA-CNN architecture. The first row indicates using some of the tricks from [11] that showed an improvement when tested on Market1501 using the HA-CNN architecture. These include warm-up, random erasing, and no-bias in the fully connected layers. These tricks alone (experiment a) provided an improvement of 2% in rank-1 accuracy and 6.3% in mean average precision compared to the original HA-CNN paper result (i.e. our baseline).

Next, to test the influence of some of our modifications we report the performance after disabling them. The most significant decrease in results compared to column i was caused by disabling the weighted triplet loss and soft margin (using the original triplet loss as in equation (1) instead) with a drop of 1.6% in rank-1 accuracy and 2.8% in mAP (column b). Cancelling the  $L_2$  normalization caused a decrease of 1.2% in rank-1 accuracy and 2.4% in mAP (column c). Reduction of other modifications such as shuffle blocks, soft margin, GeM, and deeper and wider network caused a decrease in the performance as well indicating the benefit of using it.

Finally, we used the SWAG in two experiments: experiment h and the final Compact-ReID. Continuing the training with SWAG provided an improvement in both rank-1 and mAP in both experiments. The SWAG is used in this study as a post process for models that already achieve high accuracy to show its contribution on top of that.

#### 4.3 Exploring SWAG

Our empirical experiments showed that the SWAG method consistently improved our model performance. However, it requires additional training time and uses a custom made cosine annealing learning scheme with a decay factor. Therefore, we wanted to further explore the SWAG contribution by analyzing some of our experimental results. Table 4 shows the results when testing the learning rate scheme with and without SWAG for three different setups. In the first setup we used our proposed architecture minus three main

Type	Method	Market1501		DukeMTMC	
		r = 1	mAP	r = 1	mAP
Mask-guided	SPReID [9]	92.5	81.3	84.4	71.0
	MaskReID [14]	90.0	75.3	78.8	61.9
Stripe-based	AlignedReID [29]	90.6	77.7	81.2	67.4
	SCPNet [5]	91.2	75.2	80.3	62.6
	LocalCNN [28]	91.5	77.7	82.2	66.0
	Pyramid[32]	92.8	82.1	-	-
	PCB [23]	93.8	81.6	83.3	69.2
	BFE[3]	94.5	85.0	88.7	75.8
	MGN [26]	95.7	86.9	88.7	78.4
	Pyramid[32]	95.7	88.2	89.0	79.0
	LocalCNN (MG) [28]	95.9	87.4	-	-
Dense-semantics	DSA [30]	95.7	87.6	86.2	74.3
GAN-based	Camstyle [37]	88.1	68.7	75.3	53.5
	PN-GAN [15]	89.4	72.6	73.6	53.2
	DG-Net [34]	94.8	86.0	86.6	74.8
Global feature	IDE [35]	79.5	59.9	-	-
	SVDNet [22]	82.3	62.1	76.7	56.8
	TriNet[6]	84.9	69.1	-	-
	AWTL[19]	89.5	75.7	79.8	63.4
	OS-Net [38]	94.8	84.9	88.6	73.5
	BagOfTricks [11]	94.5	85.9	86.4	76.4
NAS	Auto-ReID [16]	94.5	85.1	88.5	75.1
Attention-based	HA-CNN [10]	91.2	75.7	80.5	63.8
	DuATM [21]	91.4	76.6	81.2	62.3
	Mancs [25]	93.1	82.3	84.9	71.8
	ABD [2]	95.6	88.3	89.0	78.6
	RGA-SC [31]	95.8	88.1	86.1	74.9
	<b>Ours (2.9M)</b>	95.8	88.7	88.8	78.9
<b>Ours (6.4M)</b>	<b>96.2</b>	<b>89.7</b>	<b>89.8</b>	<b>80.3</b>	

Table 2: Comparison of state-of-the-arts methods.

modifications: GeM, Shuffle blocks, and deeper and wider. The second and third setups are experiments g and i in Table 3 respectively. Evidently, adding the LR scheme provided a nice improvement, and adding the SWAG performed even better. The most significant improvements were in terms of mAP.

Figure 6 presents the average over five experiments comparing SWAG and standard SGD in terms of Rank1 accuracy and mAP on Market1501 dataset. Using SWAG the accuracy trend seems more consistent compared to standard SGD. In addition, it is significantly better in terms of mAP.

## 5 APPLICATION TO MULTI OBJECT TRACKING

Although the public datasets used in this study for person ReID are valuable for comparison between different architectures and models, we wanted to evaluate the model’s applicability by using it to improve multi target multi camera tracking. Testing the model in a real world setting such as tracking is much more challenging. A wrong ReID assignment can affect the assignment of other persons since we only compare each query image to tracks that are not active (not present in the room at the time of the query). In addition, for each query we need to decide if we open a new track or assign it to an existing track (ReID), meaning that in some cases the gallery does not include images of the person found in the query.

We used the LAB sequence which is a part of the Task-Decomposition database [7] of multi-view sequences for people tracking. The LAB

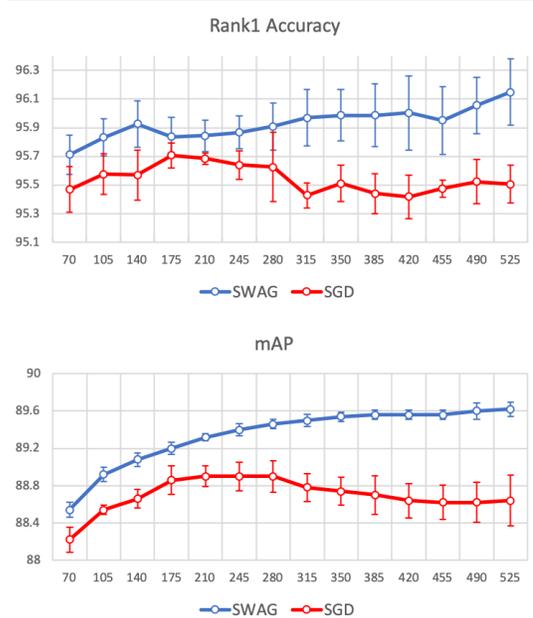


Figure 6: Performance evaluation of SWAG compared to SGD using the cosine annealing learning scheme on Market1501 dataset showing the average of 5 runs. Top: rank-1 accuracy vs. epoch. Bottom: mAP vs. epoch.

sequence is about 12.5 minutes long<sup>1</sup>, the tracking domain is about 5\*6 meters in dimension, and the images were captured at 15 Hz with a resolution of 640\*480 pixels, where four cameras are installed at the corners of the room. Through the sequence, people enter, walk around, sit down and exit the room randomly, causing frequent occlusions. The maximum number of people in the scene at the same time is 7. We first used an internal software for global people tracking which uses the calibration provided for each camera and report the results we got with and without using the model for ReID in terms of MOTA and IDF1. We used ReID each time a person enters the room by comparing it to several images per person that is currently not tracked inside the room.

Table 5 shows the results obtained using different models including: the original HA-CNN and our proposed model. Our model performed better than the original HA-CNN in terms of IDF1 using Market1501 or DukeMTMC for training. Due to the original resolution of the videos, the size of the bounding box of each query and gallery image can get very small in size. Our model showed robustness to the low-res images since it was trained on small sized input.

## 6 CONCLUSIONS

This paper explores several training techniques and architecture modifications focusing on a small-sized randomly initialized attention network for person ReID. Each training technique is tested as well

<sup>1</sup>Information in the database website mentions 3.5 minutes but the downloaded videos are actually 12.5 minutes long.

	HA-CNN [10]	a	b	c	d	e	f	g	h	i	Compact-ReID
BagOfTricks [11]		✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Soft triplet				✓	✓	✓	✓	✓	✓	✓	✓
L2 normalization			✓		✓	✓	✓	✓	✓	✓	✓
Shuffle blocks			✓	✓		✓	✓	✓	✓	✓	✓
Soft margin				✓	✓	✓	✓	✓	✓	✓	✓
GeM			✓	✓	✓	✓		✓	✓	✓	✓
Deeper & wider			✓	✓	✓	✓	✓			✓	✓
SWAG									✓		✓
Market1501 Rank1	91.2	93.2	94.1	94.5	94.9	95.1	95.3	95.4	95.8	95.7	96.2
Market1501 mAP	75.7	82.0	85.3	85.7	86.6	87.1	87.9	87.3	88.7	88.1	89.7

**Table 3: Ablation study on Market1501.** The first column indicates the different training techniques and architecture modifications we tried including some of the tricks mentioned in BagOfTricks [11]: warmup, random erase, label smoothing, and no bias in the classification layers. The baseline we started with, i.e. the original HA-CNN implementation, is presented in the second column for comparison. The last column shows the results of our proposed Compact-ReID network including all of the training techniques and architecture modifications proposed in this study. Columns a-i demonstrates the impact of each modification by turning it off.

Setup	Update	Market1501 r = 1	mAP
1	-	93.8	83.6
	+LR Scheme	94.3	84.8
	+SWAG	<b>94.5</b>	<b>85.3</b>
2	-	95.4	87.3
	+LR Scheme	<b>95.8</b>	88.2
	+SWAG	<b>95.8</b>	<b>88.7</b>
3	-	95.7	88.1
	+LR Scheme	95.7	88.9
	+SWAG	<b>96.2</b>	<b>89.7</b>

**Table 4: Performance evaluation on Market1501 for SWAG with cosine annealing with decay factor learning scheme.**

Model	Trained Dataset	MOTA	IDF1
Compact-ReID	DukeMTMC	96.1	<b>89.1</b>
Compact-ReID	Market1501	96.1	79.6
HA-CNN	DukeMTMC	96.1	78.9
HA-CNN	Market1501	96.1	65.7
No ReID	-	96.1	57.1

**Table 5: Multi camera multi target tracking results on LAB dataset using our proposed Compact-ReID model compared to the original HA-CNN.**

as some of the tricks presented in other prior works. Using the proposed training scheme and network modifications we were able to outperform SotA works achieving 96.2% rank1 accuracy and 89.7% mAP on Market1501 and 89.8% rank1 accuracy and 80.3% mAP on DukeMTMC with only 6.4M parameters. In addition, we show that even for a smaller version (2.9M parameters) we achieve state of the art results. Finally, we show the applicability of our proposed model by utilizing it to improve existing methods for multi object tracking on a public dataset. Future work entails more experiments using other deep ReID networks as our baseline, as well as tackling the cross-domain challenges in person ReID.

## ACKNOWLEDGMENTS

We would like to thank Sagi Rorlich and Genadiy Vasserman for their help in some of the experiments.

## REFERENCES

- [1] Dapeng Chen, Dan Xu, Hongsheng Li, Nicu Sebe, and Xiaogang Wang. 2018. Group consistent similarity learning via deep crf for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8649–8658.
- [2] Tianlong Chen, Shaojin Ding, Jingyi Xie, Ye Yuan, Wuyang Chen, Yang Yang, Zhou Ren, and Zhangyang Wang. 2019. ABD-Net: Attentive but Diverse Person Re-Identification. *arXiv preprint arXiv:1908.01114* (2019).
- [3] Zuozhuo Dai, Mingqiang Chen, Siyu Zhu, and Ping Tan. 2018. Batch feature erasing for person re-identification and beyond. *arXiv preprint arXiv:1811.07130* (2018).
- [4] Xing Fan, Wei Jiang, Hao Luo, and Mengjuan Fei. 2019. Spherereid: Deep hypersphere manifold embedding for person re-identification. *Journal of Visual Communication and Image Representation* 60 (2019), 51–58.
- [5] Xing Fan, Hao Luo, Xuan Zhang, Lingxiao He, Chi Zhang, and Wei Jiang. 2018. Scpnet: Spatial-channel parallelism network for joint holistic and partial person re-identification. In *Asian Conference on Computer Vision*. Springer, 19–34.
- [6] Alexander Hermans, Lucas Beyer, and Bastian Leibe. 2017. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737* (2017).
- [7] Tao Hu, Stefano Messelodi, and Oswald Lanz. 2014. Dynamic task decomposition for probabilistic tracking in complex scenes. In *2014 22nd International Conference on Pattern Recognition*. IEEE, 4134–4139.
- [8] Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. 2018. Averaging weights leads to wider optima and better generalization. *arXiv preprint arXiv:1803.05407* (2018).
- [9] Mahdi M Kalayeh, Emrah Basaran, Muhittin Gökmen, Mustafa E Kamasak, and Mubarak Shah. 2018. Human semantic parsing for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1062–1071.
- [10] Wei Li, Xiatian Zhu, and Shaogang Gong. 2018. Harmonious attention network for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2285–2294.
- [11] Hao Luo, Youzhi Gu, Xingyu Liao, Shenqi Lai, and Wei Jiang. 2019. Bag of tricks and a strong baseline for deep person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. 0–0.
- [12] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. 2018. ShuffleNet v2: Practical guidelines for efficient cnn architecture design. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 116–131.
- [13] Wesley Maddox, Timur Garipov, Pavel Izmailov, Dmitry Vetrov, and Andrew Gordon Wilson. 2019. A simple baseline for bayesian uncertainty in deep learning. *arXiv preprint arXiv:1902.02476* (2019).
- [14] Lei Qi, Jing Huo, Lei Wang, Yinghuan Shi, and Yang Gao. 2018. Maskreid: A mask based deep ranking neural network for person re-identification. *arXiv preprint arXiv:1804.03864* (2018).
- [15] Xuelin Qian, Yanwei Fu, Tao Xiang, Wenxuan Wang, Jie Qiu, Yang Wu, Yu-Gang Jiang, and Xiangyang Xue. 2018. Pose-normalized image generation for person

- re-identification. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 650–667.
- [16] Ruijie Quan, Xuanyi Dong, Yu Wu, Linchao Zhu, and Yi Yang. 2019. Auto-ReID: Searching for a Part-aware ConvNet for Person Re-Identification. *arXiv preprint arXiv:1903.09776* (2019).
- [17] Filip Radenović, Giorgos Tolias, and Ondřej Chum. 2018. Fine-tuning CNN image retrieval with no human annotation. *IEEE transactions on pattern analysis and machine intelligence* 41, 7 (2018), 1655–1668.
- [18] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi. 2016. Performance measures and a data set for multi-target, multi-camera tracking. In *European Conference on Computer Vision*. Springer, 17–35.
- [19] Ergys Ristani and Carlo Tomasi. 2018. Features for multi-target multi-camera tracking and re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 6036–6046.
- [20] Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 815–823.
- [21] Jianlou Si, Honggang Zhang, Chun-Guang Li, Jason Kuen, Xiangfei Kong, Alex C Kot, and Gang Wang. 2018. Dual attention matching network for context-aware feature sequence based person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5363–5372.
- [22] Yifan Sun, Liang Zheng, Weijian Deng, and Shengjin Wang. 2017. Svdnet for pedestrian retrieval. In *Proceedings of the IEEE International Conference on Computer Vision*. 3800–3808.
- [23] Yifan Sun, Liang Zheng, Yi Yang, Qi Tian, and Shengjin Wang. 2018. Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline). In *Proceedings of the European Conference on Computer Vision (ECCV)*. 480–496.
- [24] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2818–2826.
- [25] Cheng Wang, Qian Zhang, Chang Huang, Wenyu Liu, and Xinggang Wang. 2018. Mancs: A multi-task attentional network with curriculum sampling for person re-identification. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 365–381.
- [26] Guanshuo Wang, Yufeng Yuan, Xiong Chen, Jiwei Li, and Xi Zhou. 2018. Learning discriminative features with multiple granularities for person re-identification. In *2018 ACM Multimedia Conference on Multimedia Conference*. ACM, 274–282.
- [27] Di Wu, Si-Jia Zheng, Xiao-Ping Zhang, Chang-An Yuan, Fei Cheng, Yang Zhao, Yong-Jun Lin, Zhong-Qiu Zhao, Yong-Li Jiang, and De-Shuang Huang. 2019. Deep learning-based methods for person re-identification: A comprehensive review. *Neurocomputing* (2019).
- [28] Jiwei Yang, Xu Shen, Xinmei Tian, Houqiang Li, Jianqiang Huang, and Xian-Sheng Hua. 2018. Local convolutional neural networks for person re-identification. In *2018 ACM Multimedia Conference on Multimedia Conference*. ACM, 1074–1082.
- [29] Xuan Zhang, Hao Luo, Xing Fan, Weilai Xiang, Yixiao Sun, Qiqi Xiao, Wei Jiang, Chi Zhang, and Jian Sun. 2017. Alignedreid: Surpassing human-level performance in person re-identification. *arXiv preprint arXiv:1711.08184* (2017).
- [30] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, and Zhibo Chen. 2019. Densely Semantically Aligned Person Re-Identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 667–676.
- [31] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, Xin Jin, and Zhibo Chen. 2019. Relation-Aware Global Attention. *arXiv preprint arXiv:1904.02998* (2019).
- [32] Feng Zheng, Xing Sun, Xinyang Jiang, Xiaowei Guo, Zongqiao Yu, and Feiyue Huang. 2018. A coarse-to-fine pyramidal model for person re-identification via multi-loss dynamic training. *arXiv preprint arXiv:1810.12193* (2018).
- [33] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. 2015. Scalable person re-identification: A benchmark. In *Proceedings of the IEEE international conference on computer vision*. 1116–1124.
- [34] Zhedong Zheng, Xiaodong Yang, Zhiding Yu, Liang Zheng, Yi Yang, and Jan Kautz. 2019. Joint discriminative and generative learning for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2138–2147.
- [35] Zhedong Zheng, Liang Zheng, and Yi Yang. 2018. A discriminatively learned CNN embedding for person reidentification. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)* 14, 1 (2018), 13.
- [36] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. 2017. Random erasing data augmentation. *arXiv preprint arXiv:1708.04896* (2017).
- [37] Zhun Zhong, Liang Zheng, Zhedong Zheng, Shaozi Li, and Yi Yang. 2018. Camstyle: A novel data augmentation method for person re-identification. *IEEE Transactions on Image Processing* 28, 3 (2018), 1176–1190.
- [38] Kaiyang Zhou, Yongxin Yang, Andrea Cavallaro, and Tao Xiang. 2019. Omni-Scale Feature Learning for Person Re-Identification. *arXiv preprint arXiv:1905.00953* (2019).