# In Defense of the Triplet Loss for Person Re-Identification

Alexander Hermans,* Lucas Beyer* and Bastian Leibe
Visual Computing Institute
RWTH Aachen University
last@vision.rwth-aachen.org

## Abstract

*In the past few years, the field of computer vision has gone through a revolution fueled mainly by the advent of large datasets and the adoption of deep convolutional neural networks for end-to-end learning. The person re-identification subfield is no exception to this. Unfortunately, a prevailing belief in the community seems to be that the triplet loss is inferior to using surrogate losses (classification, verification) followed by a separate metric learning step. We show that, for models trained from scratch as well as pretrained ones, using a variant of the triplet loss to perform end-to-end deep metric learning outperforms most other published methods by a large margin.*

## 1. Introduction

In recent years, person re-identification (ReID) has attracted significant attention in the computer vision community. Especially with the rise of deep learning, many new approaches have been proposed to achieve this task [40, 8, 42, 31, 39, 10, 52, 4, 46, 20, 54, 35][1]. In many aspects person ReID is similar to image retrieval, where significant progress has been made and where deep learning has recently introduced a lot of changes. One prominent example in the recent literature is FaceNet [29], a convolutional neural network (CNN) used to learn an embedding for faces. The key component of FaceNet is to use the triplet loss, as introduced by Weinberger and Saul [41], for training the CNN as an embedding function. The triplet loss optimizes the embedding space such that data points with the same identity are closer to each other than those with different identities. A visualization of such an embedding is shown in Figure 1.

Several approaches for person ReID have already used some variant of the triplet loss to train their models [17, 9, 28, 8, 40, 31, 33, 26, 6, 25], with moderate success. The

---

*Equal contribution. Ordering determined by a last minute coin flip.
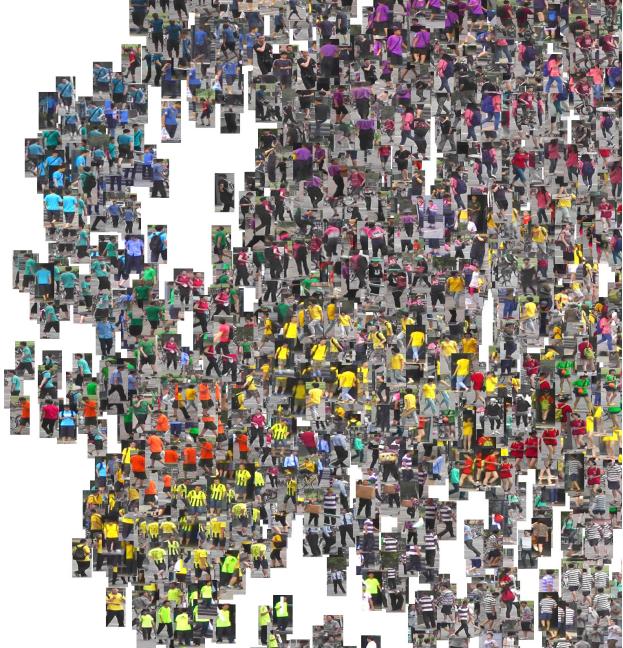[1]A nice overview of the field is given by a recent survey paper [51].



Figure 1: A small crop of the Barnes-Hut t-SNE [38] of our learned embeddings for the Market-1501 test-set. The triplet loss learns semantically meaningful features.

recently most successful person ReID approaches argue that a classification loss, possibly combined with a verification loss, is superior for the task [6, 51, 10, 52, 22]. Typically, these approaches train a deep CNN using one or multiple of these surrogate losses and subsequently use a part of the network as a feature extractor, combining it with a metric learning approach to generate final embeddings. Both of these losses have their problems, though. The classification loss necessitates a growing number of learnable parameters as the number of identities increases, most of which will be discarded after training. On the other hand, many of the networks trained with a verification loss have to be used in a cross-image representation mode, only answering the question "How similar are these two images?". This makes

using them for any other task, such as clustering or retrieval, prohibitively expensive, as each probe has to go through the network paired up with every gallery image.

In this paper we show that, contrary to current opinion, a plain CNN with a triplet loss can outperform current state-of-the-art approaches on the CUHK03 [21], Market-1501 [50] and MARS [49] datasets. The triplet loss allows us to perform end-to-end learning between the input image and the desired embedding space. This means we directly optimize the network for the final task, which renders an additional metric learning step obsolete. Instead, we can simply compare persons by computing the Euclidean distance of their embeddings.

A possible reason for the unpopularity of the triplet loss is that, when applied naïvely, it will indeed often produce disappointing results. An essential part of learning using the triplet loss is the mining of hard triplets, as otherwise training will quickly stagnate. However, mining such hard triplets is time consuming and it is unclear what defines "good" hard triplets [29, 31]. Even worse, selecting too hard triplets too often makes the training unstable. We show how this problem can be alleviated, resulting in both faster training and better performance. We systematically analyze the design space of triplet losses, and evaluate which one works best for person ReID. While doing so, we place two previously proposed variants [9, 32] into this design space and discuss them in more detail in Section 2. Specifically, we find that the best performing version has not been used before. Furthermore we also show that a margin-less formulation performs slightly better, while removing one hyper-parameter.

Another clear trend seems to be the use of pretrained models such as GoogleNet [36] or ResNet-50 [14]. Indeed, pretrained models often obtain great scores for person ReID [10, 52], while ever fewer top-performing approaches use networks trained from scratch [21, 1, 8, 42, 31, 39, 4]. Some authors even argue that training from scratch is bad [10]. However, using pretrained networks also leads to a design lock-in, and does not allow for the exploration of new deep learning advances or different architectures. We show that, when following best practices in deep learning, networks trained from scratch can perform competitively for person ReID. Furthermore, we do not rely on network components specifically tailored towards person ReID, but train a plain feed-forward CNN, unlike many other approaches that train from scratch [21, 1, 39, 42, 34, 20, 48]. Indeed, our networks using pretrained weights obtain the best results, but our far smaller architecture obtains respectable scores, providing a viable alternative for applications where person ReID needs to be performed on resource-constrained hardware, such as embedded devices.

In summary our contribution is twofold: Firstly we introduce variants of the classic triplet loss which render mining of hard triplets unnecessary and we systematically evaluate these variants. And secondly, we show how, contrary to the prevailing opinion, using a triplet loss and no special layers, we achieve state-of-the-art results both with a pretrained CNN and with a model trained from scratch. This highlights that a well designed triplet loss has a significant impact on the result, on par with other architectural novelties, hopefully enabling other researchers to gain the full potential of the previously often dismissed triplet loss. This is an important result, highlighting that a well designed triplet loss has a significant impact on model performance—on par with other architectural novelties—hopefully enabling other researchers to gain the full potential of the previously often dismissed triplet loss.

## 2. Learning Metric Embeddings, the Triplet Loss, and the Importance of Mining

The goal of metric embedding learning is to learn a function $f_\theta(x) : \mathbb{R}^F \to \mathbb{R}^D$ which maps semantically similar points from the data manifold in $\mathbb{R}^F$ onto metrically close points in $\mathbb{R}^D$. Analogously, $f_\theta$ should map semantically different points in $\mathbb{R}^F$ onto metrically distant points in $\mathbb{R}^D$. The function $f_\theta$ is parametrized by $\theta$ and can be anything ranging from a linear transform [41, 23, 45, 28] to complex non-linear mappings usually represented by deep neural networks [9, 8, 10]. Let $D(x, y) : \mathbb{R}^D \times \mathbb{R}^D \to \mathbb{R}$ be a metric function measuring distances in the embedding space. For clarity we use the shortcut notation $D_{i,j} = D(f_\theta(x_i), f_\theta(x_j))$, where we omit the indirect dependence of $D_{i,j}$ on the parameters $\theta$. As is common practice, all loss-terms are divided by the number of summands in a batch; we omit this term in the following equations for conciseness.

Weinberger and Saul [41] explore this topic with the explicit goal of performing $k$-nearest neighbor classification in the learned embedding space and propose the "Large Margin Nearest Neighbor loss" for optimizing $f_\theta$:

$$\mathcal{L}_{\text{LMNN}}(\theta) = (1 - \mu)\mathcal{L}_{\text{pull}}(\theta) + \mu\mathcal{L}_{\text{push}}(\theta), \quad (1)$$

which is comprised of a *pull*-term, pulling data points $i$ towards their *target neighbor* $T(i)$ from the same class, and a *push*-term, pushing data points from a different class $k$ further away:

$$\mathcal{L}_{\text{pull}}(\theta) = \sum_{i,j \in T(i)} D_{i,j}, \quad (2)$$

$$\mathcal{L}_{\text{push}}(\theta) = \sum_{\substack{a,n \\ y_a \neq y_n}} \left[ m + D_{a,T(a)} - D_{a,n} \right]_+. \quad (3)$$

Because the motivation was nearest-neighbor classification, allowing disparate clusters of the same class was an explicit

goal, achieved by choosing fixed target neighbors at the onset of training. Since this property is harmful for retrieval tasks such as face and person ReID, FaceNet [29] proposed a modification of $\mathcal{L}_{\text{LMNN}}(\theta)$ called the "Triplet loss":

$$\mathcal{L}_{\text{tri}}(\theta) = \sum_{\substack{a,p,n \\ y_a = y_p \neq y_n}} [m + D_{a,p} - D_{a,n}]_+ . \quad (4)$$

This loss makes sure that, given an *anchor* point $x_a$, the projection of a *positive* point $x_p$ belonging to the same class (person) $y_a$ is closer to the anchor's projection than that of a *negative* point belonging to another class $y_n$, by at least a margin $m$. If this loss is optimized over the whole dataset for long enough, eventually all possible pairs $(x_a, x_p)$ will be seen and be pulled together, making the pull-term redundant. The advantage of this formulation is that, while eventually all points of the same class will form a single cluster, they are not required to collapse to a single point; they merely need to be closer to each other than to any point from a different class.

A major caveat of the triplet loss, though, is that as the dataset gets larger, the possible number of triplets grows cubically, rendering a long enough training impractical. To make matters worse, $f_\theta$ relatively quickly learns to correctly map most trivial triplets, rendering a large fraction of all triplets uninformative. Thus mining *hard* triplets becomes crucial for learning. Intuitively, being told over and over again that people with differently colored clothes are different persons does not teach one anything, whereas seeing similarly-looking but different people (*hard negatives*), or pictures of the same person in wildly different poses (*hard positives*) dramatically helps understanding the concept of "same person". On the other hand, being shown *only* the hardest triplets would select outliers in the data unproportionally often and make $f_\theta$ unable to learn "normal" associations, as will be shown in Table 1. Examples of typical hard positives, hard negatives, and outliers are shown in the Supplementary Material. Hence it is common to only mine *moderate* negatives [29] and/or *moderate* positives [31]. Regardless of which type of mining is being done, it is a separate step from training and adds considerable overhead, as it requires embedding a large fraction of the data with the most recent $f_\theta$ and computing all pairwise distances between those data points.

In a classical implementation, once a certain set of $B$ triplets has been chosen, their images are stacked into a batch of size $3B$, for which the $3B$ embeddings are computed, which are in turn used to create $B$ terms contributing to the loss. Given the fact that there are up to $6B^2 - 4B$ possible combinations of these $3B$ images that are valid triplets, using only $B$ of them seems wasteful. With this realization, we propose an organizational modification to the classic way of using the triplet loss: the core idea is to

form batches by randomly sampling $P$ classes (person identities), and then randomly sampling $K$ images of each class (person), thus resulting in a batch of $PK$ images.[2] Now, for each sample $a$ in the batch, we can select the hardest positive and the hardest negative samples *within the batch* when forming the triplets for computing the loss, which we call *Batch Hard*:

$$\mathcal{L}_{\text{BH}}(\theta; X) = \overbrace{\sum_{i=1}^{P} \sum_{a=1}^{K}}^{\text{all anchors}} \Big[ m + \overbrace{\max_{p=1\ldots K} D\left(f_\theta(x_a^i), f_\theta(x_p^i)\right)}^{\text{hardest positive}} \quad (5)$$
$$- \underbrace{\min_{\substack{j=1\ldots P \\ n=1\ldots K \\ j \neq i}} D\left(f_\theta(x_a^i), f_\theta(x_n^j)\right)}_{\text{hardest negative}} \Big]_+ ,$$

which is defined for a mini-batch $X$ and where a data point $x_j^i$ corresponds to the $j$-th image of the $i$-th person in the batch.

This results in $PK$ terms contributing to the loss, a threefold[3] increase over the traditional formulation. Additionally, the selected triplets can be considered *moderate triplets*, since they are the hardest within a small subset of the data, which is exactly what is best for learning with the triplet loss.

This new formulation of sampling a batch immediately suggests another alternative, that is to simply use all possible $PK(PK - K)(K - 1)$ combinations of triplets, which corresponds to the strategy chosen in [9] and which we call *Batch All*:

$$\mathcal{L}_{\text{BA}}(\theta; X) = \overbrace{\sum_{i=1}^{P} \sum_{a=1}^{K}}^{\text{all anchors}} \overbrace{\sum_{\substack{p=1 \\ p \neq a}}^{K}}^{\text{all pos.}} \overbrace{\sum_{\substack{j=1 \\ j \neq i}}^{P} \sum_{n=1}^{K}}^{\text{all negatives}} \left[ m + d_{j,a,n}^{i,a,p} \right]_+ , \quad (6)$$
$$d_{j,a,n}^{i,a,p} = D\left(f_\theta(x_a^i), f_\theta(x_p^i)\right) - D\left(f_\theta(x_a^i), f_\theta(x_n^j)\right) .$$

At this point, it is important to note that both $\mathcal{L}_{\text{BH}}$ and $\mathcal{L}_{\text{BA}}$ still exactly correspond to the standard triplet loss in the limit of infinite training. Both the max and min functions are continuous and differentiable almost everywhere, meaning they can be used in a model trained by stochastic (sub-)gradient descent without concern. In fact, they are already widely available in popular deep-learning frameworks for the implementation of max-pooling and the ReLU [11] non-linearity.

Most similar to our *batch hard* and *batch all* losses is the *Lifted Embedding* loss [32], which fills the batch with

---

[2]In all experiments we choose $B$, $P$, and $K$ in such a way that $3B$ is close to $PK$, *e.g.* $3 \cdot 42 \approx 32 \cdot 4$.

[3]Because $PK \approx 3B$, see footnote 2

triplets but considers all but the anchor-positive pair as negatives:

$$\mathcal{L}_{\text{L}}(\theta; X) = \sum_{(a,p) \in X} \left[ D_{a,p} + \log \sum_{\substack{n \in X \\ n \neq a, n \neq p}} \left( e^{m - D_{a,n}} + e^{m - D_{p,n}} \right) \right]_+ .$$

While [32] motivates a "hard"-margin loss similar to $\mathcal{L}_{\text{BH}}$ and $\mathcal{L}_{\text{BA}}$, they end up optimizing the smooth bound of it given in the above equation. Additionally, traditional $3B$ batches are considered, thus using all possible negatives, but only one positive pair per triplet. This leads us to propose a generalization of the Lifted Embedding loss based on $PK$ batches which considers all anchor-positive pairs as follows:

$$\mathcal{L}_{\text{LG}}(\theta; X) = \sum_{i=1}^{P} \sum_{a=1}^{K} \left[ \log \sum_{\substack{p=1 \\ p \neq a}}^{K} e^{D\left(f_\theta(x_a^i), f_\theta(x_p^i)\right)} \right. \qquad (7)$$

$$\left. + \log \sum_{\substack{j=1 \\ j \neq i}}^{P} \sum_{n=1}^{K} e^{m - D\left(f_\theta(x_a^i), f_\theta(x_n^j)\right)} \right]_+ .$$

**Distance Measure.** Throughout this section, we have referred to $D(a, b)$ as the distance function between $a$ and $b$ in the embedding space. In most related works, the squared Euclidean distance $D\left(f_\theta(x_i), f_\theta(x_j)\right) = \|f_\theta(x_i) - f_\theta(x_j)\|_2^2$ is used as metric, although nothing in the above loss definitions precludes using any other (sub-)differentiable distance measure. While we do not have a side-by-side comparison, we noticed during initial experiments that using the squared Euclidean distance made the optimization more prone to collapsing, whereas using the actual (non-squared) Euclidean distance was more stable. We hence used the Euclidean distance throughout all our experiments presented in this paper. In addition, squaring the Euclidean distance makes the margin parameter less interpretable, as it does not represent an absolute distance anymore.

Note that when forcing the embedding's norm to one, using the squared Euclidean distance corresponds to using the cosine-similarity, up to a constant factor of two. We did not use a normalizing layer in any of our final experiments. For one, it does not dramatically regularize the network by reducing the available embedding space: the space spanned by all $D$-dimensional vector of fixed norm is still a $D - 1$-dimensional volume. Worse, an output-normalization layer can actually hide problems in the training, such as slowly collapsing or exploding embeddings.

**Soft-margin.** The role of the hinge function $[m + \bullet]_+$ is to avoid correcting "already correct" triplets. But in person ReID, it can be beneficial to pull together samples from the same class as much as possible [45, 8], especially when working on tracklets such as in MARS [49]. For this purpose, it is possible to replace the hinge function by a smooth approximation using the softplus function: $\ln(1 + \exp(\bullet))$, for which numerically stable implementations are commonly available as `log1p`. The softplus function has similar behavior to the hinge, but it decays exponentially instead of having a hard cut-off, we hence refer to it as the soft-margin formulation.

**Summary.** In summary, the novel contributions proposed in this paper are the *batch hard* loss and its soft margin version. In the following section we evaluate them experimentally and show that, for ReID, they achieve superior performance compared to both the traditional triplet loss and the previously published variants of it [9, 32].

## 3. Experiments

Our experimental evaluation is split up into three main parts. The first section evaluates different variations of the triplet loss, including some hyper-parameters, and identifies the setting that works best for person ReID. This evaluation is performed on a train/validation split we create based on the MARS training set. The second section shows the performance we can attain based on the selected variant of the triplet loss. We show state-of-the-art results on the CUHK03, Market-1501 and MARS test sets, based on a pretrained network and a network trained from scratch. Finally, the third section discusses advantages of training models from scratch with respect to real-world use cases.

### 3.1. Datasets

We focus on the Market-1501 [50] and MARS [49] datasets, the two largest person ReID datasets currently available. The Market-1501 dataset contains bounding boxes from a person detector which have been selected based on their intersection-over-union overlap with manually annotated bounding boxes. It contains $32\,668$ images of $1501$ persons, split into train/test sets of $12\,936/19\,732$ images as defined by [50]. The dataset uses both single- and multi-query evaluation, we report numbers for both.

The MARS dataset originates from the same raw data as the Market-1501 dataset; however, a significant difference is that the MARS dataset does not have any manually annotated bounding boxes, reducing the annotation overhead. MARS consist of "tracklets" which have been grouped into person IDs manually. It contains $1\,191\,003$ images split into train/test sets of $509\,914/681\,089$ images, as defined by [49]. Here, person ReID is no longer performed on a frame-to-frame level, but instead on a tracklet-to-tracklet level, where feature embeddings are pooled across a tracklet, thus it is inherently a multi-query setup.

We use the standard evaluation metrics for both datasets, namely the mean average precision score (mAP) and the cumulative matching curve (CMC) at rank-1 and rank-5. To compute these scores we use the evaluation code provided by [55].

Additionally, we show results on the CUHK03 [21] dataset for our pretrained network, using the single shot setup and average over the provided 20 train/test splits.

## 3.2. Training

Unless specifically noted otherwise, we use the same training procedure across all experiments and on all datasets. We performed all our experiments using the Theano [5] framework, code is available at `redacted`.

We use the Adam optimizer [18] with the default hyper-parameter values ($\epsilon = 10^{-3}$, $\beta_1 = 0.9$, $\beta_2 = 0.999$) for most experiments. During initial experiments on our own MARS validation split (see Sec. 3.4), we ran multiple experiments for a very long time and monitored the loss and mAP curves. With this information, we decided to fix the following exponentially decaying training schedule, which does not disadvantage any setup, for all experiments presented in this paper:

$$\epsilon(t) = \begin{cases} \epsilon_0 & \text{if } t \leq t_0 \\ \epsilon_0 0.001^{\frac{t-t_0}{t_1-t_0}} & \text{if } t_0 \leq t \leq t_1 \end{cases} \quad (8)$$

with $\epsilon_0 = 10^{-3}$, $t_0 = 15\,000$, and $t_1 = 25\,000$, stopping training when reaching $t_1$. We also set $\beta_1 = 0.5$ when entering the decay schedule at $t_0$, as is common practice [2].

In the Supplementary Material, we provide a detailed discussion of various interesting effects we regularly observed during training, providing hands-on guidance for other researchers.

## 3.3. Network Architectures

For our main results we use two different architectures, one based on a pretrained network and one which we train from scratch.

**Pretrained.** We use the ResNet-50 architecture and the weights provided by He *et al.* [14]. We discard the last layer and add two fully connected layers for our task. The first has 1024 units, followed by batch normalization [16] and ReLU [11], the second goes down to 128 units, our final embedding dimension. Trained with our *batch hard* triplet loss, we call this model *TriNet*. Due to the size of this network (25.74 M parameters), we had to limit our batch size to 72, containing $P = 18$ persons with $K = 4$ images each. For these pretrained experiments, $\epsilon_0 = 10^{-3}$ proved to be too high, causing the models to diverge within few iterations. We thus reduced $\epsilon_0$ to $3 \cdot 10^{-4}$ which worked fine on all datasets.

**Trained from Scratch.** To show that training from scratch does not necessarily result in poor performance, we also designed a network called *LuNet* which we train from scratch. *LuNet* follows the style of ResNet-v2, but uses leaky ReLU nonlinearities, multiple $3 \times 3$ max-poolings with stride 2 instead of strided convolutions, and omits the final average-pooling of feature-maps in favor of a channel-reducing final res-block. An in-depth description of the architecture is given in the Supplementary Material. As the network is much more lightweight (5.00 M parameters) than its pretrained sibling, we sample batches of size 128, containing $P = 32$ persons with $K = 4$ images each.

## 3.4. Triplet Loss

Our initial experiments test the different variants of triplet training that we discussed in Sec. 2. In order not to perform model-selection on the test set, we randomly sample a validation set of 150 persons from the MARS training set, leaving the remaining 475 persons for training. In order to make this exploration tractable, we run all of these experiments using the smaller *LuNet* trained from scratch on images downscaled by a factor of two. Since our goal here is to explore triplet loss formulations, as opposed to reaching top performance, we do not perform any data augmentation in these experiments.

Table 1 shows the resulting mAP and rank-1 scores for the different formulations at multiple margin values, and with a soft-margin where applicable. Consistent with results reported in several recent papers [10, 9, 6], the vanilla triplet loss with randomly sampled triplets performs poorly. When performing simple offline hard-mining (OHM), the scores sometimes increase dramatically, but the training also fails to learn useful embeddings for multiple margin values. This problem is well-known [29, 31] and has been discussed in Sec. 2. While the idea of learning embeddings using triplets is theoretically pleasing, this practical finnickyness, coupled with the considerable increase in training time due to non-parallelizable offline mining (from 7 h to 20 h in our experiments), makes learning with vanilla triplets rather unattractive.

Considering the long training times, it is nice to see that all proposed triplet re-formulations perform similarly to or better than the best OHM run. The key observation is that the (semi) hard-mining happens within the batch and thus comes at almost no additional runtime cost.

Perhaps surprisingly, the *batch hard* variant (Eq. 5) consistently outperforms the *batch all* variant (Eq. 6) previously used by several authors [9, 40]. We suspect this is due to the fact that in the latter, many of the possible triplets in a batch are zero, essentially "washing out" the few useful contributing terms during averaging. To test this hypothesis, we also ran experiments where we only average the non-zero loss terms (marked by $\neq 0$ in Table 1); this performs much better

| | margin 0.1 | | margin 0.2 | | margin 0.5 | | margin 1.0 | | soft margin | |
|---|---|---|---|---|---|---|---|---|---|---|
| | mAP | rank-1 | mAP | rank-1 | mAP | rank-1 | mAP | rank-1 | mAP | rank-1 |
| Triplet ($\mathcal{L}_{\mathrm{tri}}$) | 40.80 | 59.23 | 41.71 | 60.78 | 43.51 | 60.87 | 43.61 | 61.63 | *48.40* | 66.37 |
| Triplet ($\mathcal{L}_{\mathrm{tri}}$) + OHM | 16.6* | 36.6* | *61.40* | 82.95 | 32.0* | 57.1* | 41.45 | 59.42 | 46.63 | 65.43 |
| Batch hard ($\mathcal{L}_{\mathrm{BH}}$) | **65.09** | 83.51 | **65.27** | 84.55 | **65.12** | 83.39 | 63.78 | 82.48 | ***65.77*** | 84.69 |
| Batch hard ($\mathcal{L}_{\mathrm{BH}\neq0}$) | 63.10 | 83.04 | *64.19* | 83.42 | 63.71 | 82.29 | **64.06** | 84.50 | - | - |
| Batch all ($\mathcal{L}_{\mathrm{BA}}$) | 59.43 | 79.24 | 60.48 | 79.99 | 60.30 | 79.52 | *62.08* | 80.55 | 61.04 | 80.65 |
| Batch all ($\mathcal{L}_{\mathrm{BA}\neq0}$) | 63.29 | 83.65 | 64.31 | 83.37 | *64.41* | 83.98 | **64.06** | 82.90 | - | - |
| Lifted 3-pos. ($\mathcal{L}_{\mathrm{LG}}$) | 64.00 | 82.71 | 63.87 | 82.86 | 63.61 | 84.55 | *64.02* | 84.17 | - | - |
| Lifted 1-pos. ($\mathcal{L}_{\mathrm{L}}$) [32] | 61.95 | 81.35 | *63.68* | 81.73 | 63.01 | 82.48 | 62.28 | 82.34 | - | - |

Table 1: *LuNet* scores on our MARS validation split. The best performing loss at a given margin is bold, the best margin for a given loss is italic, and the overall best combination is highlighted in green. A * denotes runs trapped in a bad local optimum.

in the *batch all* case. Another interpretation of this modification is that it dynamically increases the weight of triplets which remain active as they get fewer.

The lifted triplet loss $\mathcal{L}_{\mathrm{L}}$ as introduced by [32] performs competitively, but is slightly worse than most other formulations. As can be seen in the table, our generalization to multiple positives (Eq. 7), which makes it more similar to the *batch all* variant of the triplet, improves upon it overall.

The best score was obtained by the soft-margin variation of the *batch hard* loss. We use this loss in all our further triplet experiments. To clarify, here we merely seek the best triplet loss variation for person ReID, but do not claim that this variant works best across all fields. For other tasks such as image retrieval or clustering, additional experiments will have to be performed.

### 3.5. Performance Evaluation

Here, we present the main experiments of this paper. We perform all following experiments using the *batch hard* variant $\mathcal{L}_{\mathrm{BH}}$ of the triplet loss and the soft margin, since this setup performed best during the exploratory phase.

**Batch Generation and Augmentation.** Since our *batch hard* triplet loss requires slightly different mini-batches, we sample random $PK$-style batches by first randomly sampling $P$ person identities uniformly without replacement. For each person, we then sample $K$ images, without replacement whenever possible, otherwise replicating images as necessary.

We follow common practice by using random crops and random horizontal flips during training [19, 1]. Specifically, we resize all images of size $H \times W$ to $1\frac{1}{8}(H \times W)$, of which we take random crops of size $H \times W$, keeping their aspect ratio intact. For all pretrained networks we set $H = 256, W = 128$ on Market-1501 and MARS and $H = 256, W = 96$ on CUHK03, whereas for the networks trained from scratch we set $H = 128, W = 64$.

We apply test-time augmentation in all our experiments. Following [19], we deterministically average over the embeddings from five crops and their flips. This typically gives an improvement of $3\%$ in the mAP score; a more detailed analysis can be found in the Supplementary Material.

**Combination of Embeddings.** For test-time augmentation, multi-query evaluation, and tracklet-based evaluation, the embeddings of multiple images need to be combined.

While the learned clusters have no reason to be Gaussian, their convex hull is trained to only contain positive points. Thus, a convex combination of multiple embeddings cannot get closer to a negative embedding than any of the original ones, which is not the case for a non-convex combination such as max-pooling. For this reason, we suggest combining triplet-based embeddings by using their mean. For example, combining tracklet-embeddings using max-pooling led to an $11.4\%$ point decrease in mAP on MARS.

**Comparison to State-of-the-Art.** Tables 2 and 3 compare our results to a set of related, top performing approaches on Market-1501 and MARS, and CUHK03, respectively. We do not include approaches which are orthogonal to ours and could be integrated in a straightforward manner, such as various re-ranking schemes, data augmentation, and regularization [3, 43, 54, 7, 47, 53, 56]. These are included in more exhaustive tables in the Supplementary Material. The different approaches are categorized into three major types: Identification models (I) that are trained to classify person IDs, Verification models (V) that learn whether an image pair represents the same person, and methods such as ours that directly learn an Embedding (E).

We present results for both our pretrained network (TriNet) and the network we trained from scratch (LuNet). As can clearly be seen, TriNet outperforms all current methods. Especially striking is the jump from $41.5\%$ mAP, obtained by another ResNet-50 model trained with triplets (DTL [10]), to our $69.14\%$ mAP score in Table 2. Since Geng *et al.* [10] do not discus all details of their training procedure when using the triplet loss, we could only speculate about the reasons for the large performance gap.

| | Type | Market-1501 SQ | | | Market-1501 MQ | | | MARS | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | mAP | rank-1 | rank-5 | mAP | rank-1 | rank-5 | mAP | rank-1 | rank-5 |
| TriNet | E | **69.14** | **84.92** | **94.21** | **76.42** | **90.53** | **96.29** | **67.70** | **79.80** | **91.36** |
| LuNet | E | 60.71 | 81.38 | 92.34 | 69.07 | 87.11 | 95.16 | 60.48 | 75.56 | 89.70 |
| IDE (R) + ML ours | I | 58.06 | 78.50 | 91.18 | 67.48 | 85.45 | 94.12 | 57.42 | 72.42 | 86.21 |
| LOMO + Null Space [45] | E | 29.87 | 55.43 | - | 46.03 | 71.56 | - | - | - | - |
| Gated siamese CNN [39] | V | 39.55 | 65.88 | - | 48.45 | 76.04 | - | - | - | - |
| CAN [25] | E | 35.9 | 60.3 | - | 47.9 | 72.1 | - | - | - | - |
| JLML [22] | I | 65.5 | 85.1 | - | 74.5 | 89.7 | - | - | - | - |
| ResNet 50 (I+V)$^\dagger$ [52] | I+V | 59.87 | 79.51 | 90.91 | 70.33 | 85.84 | 94.54 | - | - | - |
| DTL$^\dagger$ [10] | E | 41.5 | 63.3 | - | 49.7 | 72.4 | - | - | - | - |
| DTL$^\dagger$ [10] | I+V | 65.5 | 83.7 | - | 73.8 | 89.6 | - | - | - | - |
| APR (R, 751)$^\dagger$ [24] | I | 64.67 | 84.29 | 93.20 | - | - | - | - | - | - |
| Latent Parts (Fusion) [20] | I | 57.53 | 80.31 | - | 66.70 | 86.79 | - | 56.05 | 71.77 | 86.57 |
| IDE (R) + ML [55] | I | 49.05 | 73.60 | - | - | - | - | 55.12 | 70.51 | - |
| spatial temporal RNN [57] | E | - | - | - | - | - | - | 50.7 | 70.6 | 90.0 |
| CNN + Video$^\dagger$ [46] | I | - | - | - | - | - | - | - | 55.5 | 70.2 |
| TriNet (Re-ranked) | E | **81.07** | **86.67** | **93.38** | **87.18** | **91.75** | **95.78** | **77.43** | **81.21** | **90.76** |
| LuNet (Re-ranked) | E | 75.62 | 84.59 | 91.89 | 82.61 | 89.31 | 94.48 | 73.68 | 78.48 | 88.74 |
| IDE (R) + ML ours (Re-ra.) | I | 71.38 | 81.62 | 89.88 | 79.78 | 86.79 | 92.96 | 69.50 | 74.39 | 85.86 |
| IDE (R) + ML (Re-ra.) [55] | I | 63.63 | 77.11 | - | - | - | - | 68.45 | 73.94 | - |

Table 2: Scores on both the Market-1501 and MARS datasets. The top and middle contain our scores and those of the current state-of-the-art respectively. The bottom contains several methods with re-ranking [55]. The different types represent the optimization criteria, where I stands for identification, V for verification and E for embedding. All our scores include test-time augmentation. The best scores are bold. $^\dagger$: Concurrent work only published on arXiv.

Our LuNet model, which we train from scratch, also performs very competitively, matching or outperforming most other baselines. While it does not quite reach the performance of our pretrained model, our results clearly show that with proper training, the flexibility of training models from scratch (see Sec. 3.6) should not be discarded.

To show that the actual performance boost is indeed gained by the triplet loss and not by other design choices, we train a ResNet-50 model with a classification loss. This model is very similar to the one used in [55] and we thus refer to it as "IDE (R) ours", for which we also apply a metric learning step (XQDA [23]). Unfortunately, especially difficult images caused frequent spikes in the loss, which ended up harming the optimization using Adam. After unsuccessfully trying lower learning rates and clipping extreme loss values, we resorted to Adadelta [44], another competitive optimization algorithm which did not exhibit these problems. While we combine embeddings through average pooling for our triplet based models, we found max-pooling and normalization to work better for the classification baseline, consistent with results reported in [49]. As Table 2 shows, the performance of the resulting model "IDE (R) ours" is still on-par with similar models in the literature. However, the large gap between the identification-based model and our TriNet clearly demonstrates the advantages of using a triplet loss.

In line with the general trend in the vision community, all deep learning methods outperform shallow methods using hand-crafted features. While Table 2 only shows [45] as a non-deep learning method, to the best of our knowledge all others perform worse.

We also evaluated how our models fare when combined with a recent re-ranking approach by Zhong *et al.* [55]. This approach can be applied on top of any ranking methods and uses information from nearest neighbors in the gallery to improve the ranking result. As Table 2 shows, our approaches go well with this method and show similar improvements to those obtained by Zhong *et al.* [55].

Finally, we evaluate our models on Market-1501 with the provided $500k$ additional distractor images. The full experiment is described in the Supplementary Material. Even with these additional distractors, our triplet-based model outperforms a classification one by $8.4\%$ mAP.

All of these results show that triplet loss embeddings are indeed a valuable tool for person ReID and we expect them to significantly change the way how research will progress in this field.

| | Type | Labeled | | Detected | |
|---|---|---|---|---|---|
| | | r-1 | r-5 | r-1 | r-5 |
| TriNet | E | **89.63** | **99.01** | **87.58** | **98.17** |
| Gated siamese CNN [39] | V | - | - | 61.8 | 86.7 |
| LOMO + Null Space [45] | E | 62.55 | 90.05 | 54.70 | 84.75 |
| CAN [25] | E | 77.6 | 95.2 | 69.2 | 88.5 |
| Latent Parts (Fusion) [20] | I | 74.21 | 94.33 | 67.99 | 91.04 |
| Spindle Net* [48] | I | 88.5 | 97.8 | - | - |
| JLML [22] | I | 83.2 | 98.0 | 80.6 | 96.9 |
| DTL† [10] | I+V | 85.4 | - | 84.1 | - |
| ResNet 50 (I+V)† [52] | I+V | - | - | 83.4 | 97.1 |

Table 3: Scores on CUHK03 for TriNet and a set of recent top performing methods. The best scores are highlighted in bold. †: Concurrent work only published on arXiv. *: The method was trained on several additional datasets.

| | TriNet | | | LuNet | | |
|---|---|---|---|---|---|---|
| | mAP | rank-1 | rank-5 | mAP | rank-1 | rank-5 |
| $256 \times 128$ | 69.14 | 84.92 | 94.21 | - | - | - |
| $128 \times 64$ | 62.52 | 79.45 | 91.06 | 60.71 | 81.38 | 92.34 |
| $64 \times 32$ | 47.42 | 68.08 | 85.84 | 57.18 | 78.21 | 90.94 |

Table 4: The effect of input size on mAP and CMC scores.

## 3.6. To Pretrain or not to Pretrain?

As mentioned before, many methods for person ReID rely on pretrained networks, following a general trend in the computer vision community. Indeed, these models lead to impressive results, as we also confirmed in this paper with our TriNet model. However, pretrained networks reduce the flexibility to try out new advances in deep learning or to make task-specific changes in a network. Our LuNet model clearly suggests that it is also possible to train models from scratch and obtain competitive scores.

In particular, an interesting direction for ReID could be the usage of additional input channels such as depth information, readily available from cheap consumer hardware. However it is unclear how to best integrate such input data into a pretrained network in a proper way.

Furthermore, the typical pretrained networks are designed with accuracy in mind and do not focus on the memory footprint or the runtime of a method. Both are important factors for real-world robotic scenarios, where typically power consumption is a constraint and only less powerful hardware can be considered [12, 37]. When designing a network from scratch, one can directly take this into consideration and create networks with a smaller memory footprint and faster evaluation times.

In principle, our pretrained model can easily be sped up by using half or quarter size input images, since the global average pooling in the ResNet will still produce an output vector of the same shape. This, however, goes hand in hand with the question of how to best adapt a pretrained network to a new task with different image sizes. The typical way of leveraging pretrained networks is to simply stretch images to the fixed expected input size used to train the network, typically $224 \times 224$ pixels. We used $256 \times 128$ instead in order to preserve the aspect ratio of the original image. However, for the Market-1501 dataset, this meant we had to upscale the images, while if we do not confine ourselves to pretrained networks we can simply adjust our architecture to the dataset size, as we did in the LuNet architecture. However, we hypothesize that a pretrained network has an "intrinsic scale," for which the learned filters work properly and thus simply using smaller input images will result in suboptimal performance. To show this, we retrain our TriNet with $128 \times 64$ and $64 \times 32$ images. As Table 4 clearly shows, the performance drops rapidly. At the original image scale, our LuNet model can almost match the mAP score of TriNet and already outperforms it when considering CMC scores. At an even smaller image size, LuNet significantly outperforms the pretrained model. Since the LuNet performance only drops by about $\sim 3\%$, the small images still hold enough data to perform ReID, but the rather rigid pretrained weights can no longer adjust to such a data change. This shows that pretrained models are not a solution for arbitrary tasks, especially when one wants to train lightweight models for small images.

## 4. Discussion

We are not the first to use the triplet loss for person ReID. Ding *et al*. [9] and Wang *et al*. [40] use a batch generation and loss formulation which is very similar to our *batch all* formulation. Wang *et al*. [40] further combine it with a pairwise verification loss. However, in the *batch all* case, it was important for us to average only over the active triplets ($\mathcal{L}_{BA\neq0}$), which they do not mention. This, in combination with their rather small networks, might explain their relatively low scores. Cheng *et al*. [8] propose an "improved triplet loss" by introducing another pull term into the loss, penalizing large distances between positive images. This formulation is in fact very similar to the original one by Weinberger and Saul [41]. We briefly experimented with a pull term, but the additional weighting hyper-parameter was not trivial to optimize and it did not improve our results. Several authors suggest learning attributes and ReID jointly [17, 33, 24, 30], some of which integrate this into their embedding dimensions. This is an interesting research direction orthogonal to our work. Several other authors also defined losses over triplets of images [28, 31, 26], however, they use losses different from the triplet loss we defend in this paper, possibly explaining their lower scores. Finally, FaceNet [29] uses a huge batch with moderate mining, which can only be done on the CPU, whereas we advocate hard mining in a small batch, which has a similar effect

to moderate mining in a large batch, while fitting on a GPU and thus making training significantly more affordable.

## 5. Conclusion

In this paper we have shown that, contrary to the prevailing belief, the triplet loss is an excellent tool for person re-identification. We propose a variant that no longer requires offline hard negative mining at almost no additional cost. Combined with a pretrained network, we set the new state-of-the-art on three of the major ReID datasets. Furthermore, we show that training networks from scratch can lead to very competitive scores. We hope that in future work the ReID community will build on top of our results and shift more towards end-to-end learning.

## References

[1] E. Ahmed, M. Jones, and T. K. Marks. An Improved Deep Learning Architecture for Person Re-Identification. In *CVPR*, pages 3908–3916, 2015. 2, 6

[2] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer Normalization. *arXiv preprint arXiv:1607.06450*, 2016. 5

[3] S. Bai, X. Bai, and Q. Tian. Scalable person re-identification on supervised smoothed manifold. *CVPR*, 2017. 6, 15, 16

[4] I. B. Barbosa, M. Cristani, B. Caputo, A. Rognhaugen, and T. Theoharis. Looking Beyond Appearances: Synthetic Training Data for Deep CNNs in Re-identification. *arXiv preprint arXiv:1701.03153*, 2017. 1, 2

[5] F. Bastien, P. Lamblin, R. Pascanu, J. Bergstra, I. J. Goodfellow, A. Bergeron, N. Bouchard, and Y. Bengio. Theano: new features and speed improvements. In *NIPS W.*, 2012. 5

[6] W. Chen, X. Chen, J. Zhang, and K. Huang. A Multi-task Deep Network for Person Re-identification. *AAAI*, 2017. 1, 5

[7] Y. Chen, X. Zhu, and G. Shaogang. Person Re-Identification by Deep Learning Multi-Scale Representations. In *ICCV W. on Cross-domain Human Identification*, 2016. 6, 15, 16

[8] D. Cheng, Y. Gong, S. Zhou, J. Wang, and N. Zheng. Person Re-Identification by Multi-Channel Parts-Based CNN with Improved Triplet Loss Function. In *CVPR*, 2016. 1, 2, 4, 8

[9] S. Ding, L. Lin, G. Wang, and H. Chao. Deep feature learning with relative distance comparison for person re-identification. *Pattern Recognition*, 48(10):2993–3003, 2015. 1, 2, 3, 4, 5, 8

[10] M. Geng, Y. Wang, T. Xiang, and Y. Tian. Deep Transfer Learning for Person Re-identification. *arXiv preprint arXiv:1611.05244*, 2016. 1, 2, 5, 6, 7, 8, 15, 16

[11] X. Glorot, A. Bordes, and Y. Bengio. Deep Sparse Rectifier Neural Networks. In *AISTATS*, 2011. 3, 5, 16

[12] N. Hawes, C. Burbridge, F. Jovan, L. Kunze, B. Lacerda, L. Mudrová, J. Young, J. Wyatt, D. Hebesberger, T. Körtner, et al. The STRANDS Project: Long-Term Autonomy in Everyday Environments. *RAM*, 24(3):146–156, 2017. 8

[13] K. He, X. Zhang, S. Ren, and J. Sun. Delving Deep Into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification. In *ICCV*, pages 1026–1034, 2015. 16

[14] K. He, X. Zhang, S. Ren, and J. Sun. Deep Residual Learning for Image Recognition. In *CVPR*, 2016. 2, 5, 11

[15] K. He, X. Zhang, S. Ren, and J. Sun. Identity Mappings in Deep Residual Networks. In *ECCV*, 2016. 16

[16] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 5, 16

[17] S. Khamis, C.-H. Kuo, V. K. Singh, V. D. Shet, and L. S. Davis. Joint Learning for Attribute-Consistent Person Re-Identification. In *ECCV*, 2014. 1, 8

[18] D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. In *ICLR*, 2015. 5

[19] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. In *NIPS*, pages 1097–1105, 2012. 6, 11

[20] D. Li, X. Chen, Z. Zhang, and K. Huang. Learning Deep Context-aware Features over Body and Latent Parts for Person Re-identification. In *CVPR*, 2017. 1, 2, 7, 8, 15, 16

[21] W. Li, R. Zhao, T. Xiao, and X. Wang. DeepReID: Deep Filter Pairing Neural Network for Person Re-Identification. In *CVPR*, 2014. 2, 5

[22] W. Li, X. Zhu, and S. Gong. Person Re-Identification by Deep Joint Learning of Multi-Loss Classification. In *IJCAI*, 2017. 1, 7, 8, 15, 16

[23] S. Liao, Y. Hu, X. Zhu, and S. Z. Li. Person Re-identification by Local Maximal Occurrence Representation and Metric Learning. In *CVPR*, 2015. 2, 7

[24] Y. Lin, L. Zheng, Z. Zheng, Y. Wu, and Y. Yang. Improving person re-identification by attribute and identity learning. *arXiv preprint arXiv:1703.07220*, 2017. 7, 8, 15

[25] H. Liu, J. Feng, M. Qi, J. Jiang, and S. Yan. End-to-End Comparative Attention Networks for Person Re-identification. *Trans. Image Proc.*, 26(7):3492–3506, 2017. 1, 7, 8, 15, 16

[26] J. Liu, Z.-J. Zha, Q. Tian, D. Liu, T. Yao, Q. Ling, and T. Mei. Multi-Scale Triplet CNN for Person Re-Identification. In *ACMMM*, 2016. 1, 8

[27] A. L. Maas, A. Y. Hannun, and A. Y. Ng. Rectifier Nonlinearities Improve Neural Network Acoustic Models. In *ICML*, 2013. 16

[28] S. Paisitkriangkrai, C. Shen, and A. van den Hengel. Learning to rank in person re-identification with metric ensembles. In *CVPR*, 2015. 1, 2, 8

[29] F. Schroff, D. Kalenichenko, and J. Philbin. FaceNet: A Unified Embedding for Face Recognition and Clustering. In *CVPR*, 2015. 1, 2, 3, 5, 8

[30] A. Schumann and R. Stiefelhagen. Person Re-Identification by Deep Learning Attribute-Complementary Information. In *CVPR Workshops*, 2017. 8

[31] H. Shi, Y. Yang, X. Zhu, S. Liao, Z. Lei, W. Zheng, and S. Z. Li. Embedding Deep Metric for Person Re-identification: A Study Against Large Variations. In *ECCV*, 2016. 1, 2, 3, 5, 8

[32] H. O. Song, Y. Xiang, S. Jegelka, and S. Savarese. Deep Metric Learning via Lifted Structured Feature Embedding. In *CVPR*, 2016. 2, 3, 4, 6

[33] C. Su, S. Zhang, J. Xing, W. Gao, and Q. Tian. Deep Attributes Driven Multi-camera Person Re-identification. In *ECCV*, 2016. 1, 8

[34] A. Subramaniam, M. Chatterjee, and A. Mittal. Deep Neural Networks with Inexact Matching for Person Re-Identification. In *NIPS*, pages 2667–2675, 2016. 2

[35] Y. Sun, L. Zheng, W. Deng, and S. Wang. SVDNet for Pedestrian Retrieval. *ICCV*, 2017. 1, 15, 16

[36] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 2

[37] R. Triebel, K. Arras, R. Alami, L. Beyer, S. Breuers, R. Chatila, M. Chetouani, D. Cremers, V. Evers, M. Fiore, et al. SPENCER: A Socially Aware Service Robot for Passenger Guidance and Help in Busy Airports. In *Field and Service Robotics*, pages 607–622. Springer, 2016. 8

[38] L. Van Der Maaten. Accelerating t-SNE using Tree-Based Algorithms. *JMLR*, 15(1):3221–3245, 2014. 1, 17

[39] R. R. Varior, M. Haloi, and G. Wang. Gated Siamese Convolutional Neural Network Architecture for Human Re-Identification. In *ECCV*, 2016. 1, 2, 7, 8, 15, 16

[40] F. Wang, W. Zuo, L. Lin, D. Zhang, and L. Zhang. Joint Learning of Single-image and Cross-image Representations for Person Re-identification. In *CVPR*, 2016. 1, 5, 8

[41] K. Q. Weinberger and L. K. Saul. Distance Metric Learning for Large Margin Nearest Neighbor Classification. *JMLR*, 10:207–244, 2009. 1, 2, 8

[42] T. Xiao, H. Li, W. Ouyang, and X. Wang. Learning Deep Feature Representations with Domain Guided Dropout for Person Re-identification. In *CVPR*, 2016. 1, 2, 16

[43] R. Yu, Z. Zhou, S. Bai, and X. Bai. Divide and Fuse: A Re-ranking Approach for Person Re-identification. *BMVC*, 2017. 6

[44] M. D. Zeiler. ADADELTA: An Adaptive Learning Rate Method. *arXiv preprint arXiv:1212.5701*, 2012. 7

[45] L. Zhang, T. Xiang, and S. Gong. Learning a Discriminative Null Space for Person Re-identification. In *CVPR*, 2016. 2, 4, 7, 8, 15, 16

[46] W. Zhang, S. Hu, and K. Liu. Learning Compact Appearance Representation for Video-based Person Re-Identification. *arXiv preprint arXiv:1702.06294*, 2017. 1, 7, 15

[47] Y. Zhang, T. Xiang, T. M. Hospedales, and H. Lu. Deep Mutual Learning. *arXiv preprint arXiv:1706.00384*, 2017. 6, 15

[48] H. Zhao, M. Tian, S. Sun, J. Shao, J. Yan, S. Yi, X. Wang, and X. Tang. Spindle Net: Person Re-identification with Human Body Region Guided Feature Decomposition and Fusion. In *CVPR*, 2017. 2, 8, 15, 16

[49] L. Zheng, Z. Bie, Y. Sun, J. Wang, C. Su, S. Wang, and Q. Tian. MARS: A Video Benchmark for Large-Scale Person Re-Identification. In *ECCV*, 2016. 2, 4, 7, 15

[50] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian. Scalable Person Re-Identification: A Benchmark. In *ICCV*, 2015. 2, 4

[51] L. Zheng, Y. Yang, and A. G. Hauptmann. Person Re-identification: Past, Present and Future. *arXiv preprint arXiv:1610.02984*, 2016. 1

[52] Z. Zheng, L. Zheng, and Y. Yang. A Discriminatively Learned CNN Embedding for Person Re-identification. *arXiv preprint arXiv:1611.05666*, 2016. 1, 2, 7, 8, 11, 12, 15, 16

[53] Z. Zheng, L. Zheng, and Y. Yang. Pedestrian Alignment Network for Large-scale Person Re-identification. *arXiv preprint arXiv:1707.00408*, 2017. 6, 15

[54] Z. Zheng, L. Zheng, and Y. Yang. Unlabeled Samples Generated by GAN Improve the Person Re-identification Baseline in vitro. *ICCV*, 2017. 1, 6, 15, 16

[55] Z. Zhong, L. Zheng, D. Cao, and S. Li. Re-ranking Person Re-identification with $k$-reciprocal Encoding. *CVPR*, 2017. 5, 7, 15

[56] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang. Random Erasing Data Augmentation. *arXiv preprint arXiv:1708.04896*, 2017. 6

[57] Z. Zhou, Y. Huang, W. Wang, L. Wang, and T. Tan. See the Forest for the Trees: Joint Spatial and Temporal Recurrent Neural Networks for Video-based Person Re-identification. In *CVPR*, 2017. 7, 15

# Supplementary Material

## A. Test-time Augmentation

As is good practice in the deep learning community [19, 14], we perform test-time augmentation. From each image, we deterministically create five crops of size $H \times W$: four corner crops and one center crop, as well as a horizontally flipped copy of each. The embeddings of all these ten images are then averaged, resulting in the final embedding for a person. Table 5 shows how five possible settings affect our scores on the Market-1501 dataset. As expected, the worst option is to scale the original images to fit the network input (first line), as this shows the network an image type it has never seen before. This option is directly followed by not using any test-time augmentation, *i.e.* just using the central crop. Simply flipping the center crop and averaging the two resulting embeddings already gives a big boost while only being twice as expensive. The four additional corner crops typically seem to be less effective, while more expensive, but using both augmentations together gives the best results. For the networks trained with a triplet loss, we gain about 3% mAP, while the network trained for identification only gains about 2% mAP. A possible explanation is that the feature space we learn with the triplet loss could be more suited to averaging than that of a classification network.

|  | TriNet | | LuNet | | IDE (R) Ours | |
|---|---|---|---|---|---|---|
|  | mAP | rank-1 | mAP | rank-1 | mAP | rank-1 |
| original | 65.48 | 82.51 | 55.61 | 76.72 | 56.29 | 77.55 |
| center | 66.29 | 84.06 | 57.08 | 77.94 | 56.26 | 77.08 |
| center + flip | 68.32 | 84.47 | 59.31 | 79.78 | 57.80 | 78.71 |
| 5 crops | 67.86 | 84.83 | 59.00 | 79.42 | 56.73 | 77.29 |
| 5 crops + flip | 69.14 | 84.92 | 60.71 | 81.38 | 58.06 | 78.50 |

Table 5: The effect of test-time augmentation on Market-1501.

## B. Hard Positives, Hard Negatives and Outliers

Figure 2 shows several outliers in the Market-1501 and MARS datasets. Some issues in MARS are caused by the tracker-based annotations where bounding boxes sometimes span two persons and the tracker partially focuses on the wrong person. Additionally, some annotation mistakes can be found in both datasets; while some are obvious, some others are indeed very hard to spot!

In Figure 3 we show some of the most difficult queries along with their top-3 retrieved images (containing hard negatives), as well as their two hardest positives. While some mistakes are easy to spot by a human, others are indeed not trivial, such as the first row in Figure 3.



Figure 2: Some outliers in the Mars (top two rows) and Market-1501 (bottom row) datasets. The first row shows high image overlap between tracklets of two persons. The second row shows a very hard example where a person was wrongly matched across tracklets. The last row shows a simple annotation mistake.

## C. Experiments with Distractors

On top of the normal gallery set, the Market-1501 dataset provides an additional $500k$ distractors recorded at another time. In order to evaluate how such distractors affect performance, we randomly sample an increasing number of distractors and add them to the original gallery set. Here we compare to the results from Zheng *et al*. [52]. Both our models show a similar behavior to that of their ResNet-50 baseline. Surprisingly, our LuNet model starts out with a slightly better mAP score than the baseline and ends up just below it, while consistently being better when considering the rank-1 score. This might indeed suggest that the inductive bias from pretraining helps during generalization to large amounts of unseen data. Nevertheless, all models seem to suffer under the increasing gallery set in a similar manner, albeit none of them fails badly. Especially the fact that in $74.70\%$ of all single-image queries the first image out of $519\,732$ gallery images is correctly retrieved is an impressive result.

For reproducibility of the 500k distractor plot (Fig. 4), Table 6 lists the values of the plot.

Figure 3: Some of the hardest queries. The leftmost column shows the query image, followed by the top 3 retrieved images and the two ground truth matches with the highest distance to the query images, *i.e.* the hardest positives. Correctly retrieved images have a green border, mistakes (*i.e.* hard negatives) have a red border.

| Gallery size | TriNet | | LuNet | | Res50 (I+V) [52] | |
| --- | --- | --- | --- | --- | --- | --- |
| | mAP | rank-1 | mAP | rank-1 | mAP | rank-1 |
| 19 732 | 69.14 | 84.92 | 60.71 | 81.38 | 59.87 | 79.51 |
| 119 732 | 61.93 | 79.69 | 52.73 | 75.65 | 52.28 | 73.78 |
| 219 732 | 58.74 | 77.88 | 49.44 | 73.40 | 49.11 | 71.50 |
| 319 732 | 56.58 | 76.34 | 47.17 | 71.85 | NaN | NaN |
| 419 732 | 54.97 | 75.50 | 45.57 | 70.52 | NaN | NaN |
| 519 732 | 53.63 | 74.70 | 44.26 | 69.74 | 45.24 | 68.26 |

Table 6: Values for the 500k distractor plot.

## D. Notes on Network Training

Here we present and discuss several training-logs that display interesting behavior, as noted in the main paper.
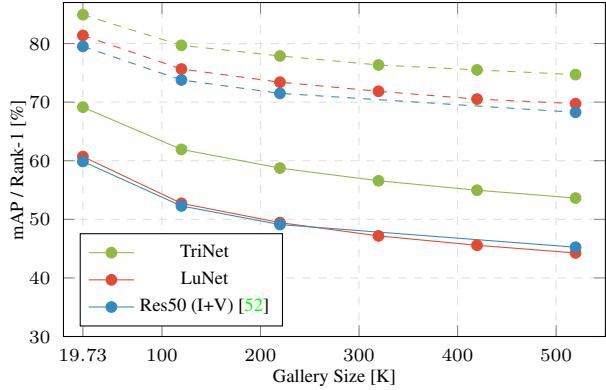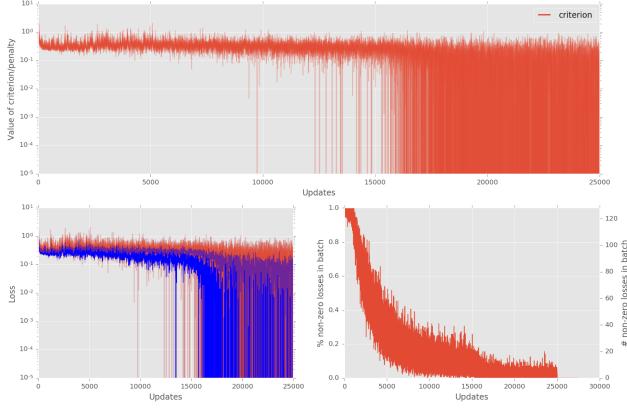


Figure 4: 500k distractor set results. Solid lines represent the mAP score, dashed lines the rank-1 score. See Supplementary Material for values.

This serves as practical advice of what to monitor for practitioners who choose to use a triplet-based loss in their training.
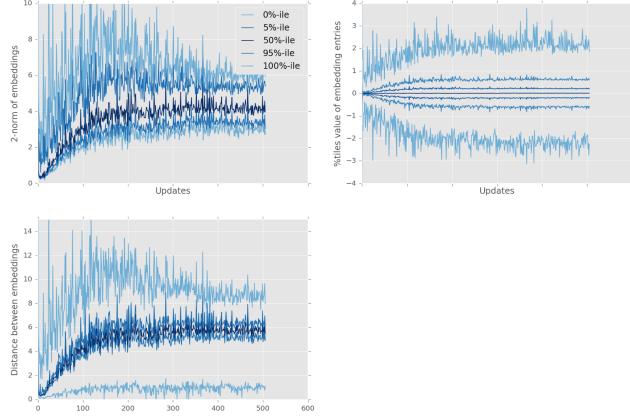
A typical training usually proceeds as follows: initially, all embeddings are pulled together towards their center of gravity. When they come close to each other, they will "pass" each other to join "their" clusters and, once this cross-over has happened, training mostly consists of pushing the clusters further apart and fine-tuning them. The *collapsing* of training happens when the margin is too large and the initial spread is too small, such that the embeddings get stuck when trying to pass each other.

Most importantly of all, if any type of hard-triplet mining is used, a stagnating loss curve by no means indicates stagnating progress. As the network learns to solve some hard cases, it will be presented with other hard cases and hence still keep a high loss. We recommend observing the fraction of active triplets in a batch, as well as the norms of the embeddings and all pairwise distances.

Figures 5, 6, 7, and 8 all show different training logs. Note that while they all share the x-axis since the number of updates was kept the same throughout the experiments, the y-axes vary for clearer visualization. First, the topmost plot in each Subfigure (a) ("Value of criterion/penalty") shows all per-triplet values of the optimization criterion (the triplet loss) encountered in each mini-batch. This is shown again in the plot below it on the left ("loss"), with an overlaid light-blue line representing the batch-mean criterion value, and an overlaid dark-blue line representing the batch's 5-percentile. To the right of it, the "% non-zero losses in batch" shows how many entries in the mini-batch had non-zero loss; values are computed up to a precision of $10^{-5}$, which explains how it can be below $100\%$ in the soft-margin case. The single final vertical line present in some plots should be ignored as a plotting-artifact.

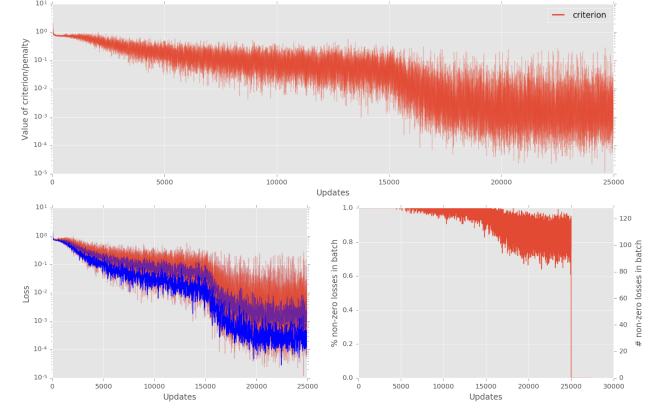(a) Training-log of the loss and active triplet count.
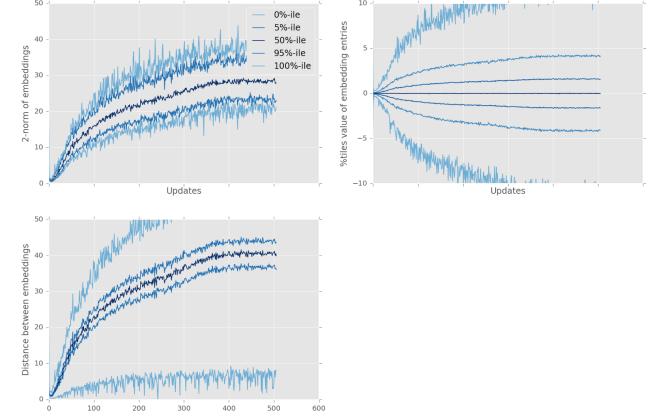


(b) Training-log of the embeddings in the minibatch.

Figure 5: Training-log of LuNet on Market1501 using the *batch hard* triplet loss with margin $0.2$. The embeddings stay bounded, as expected from a triplet formulation, and there is a lot of progress even when the loss stays seemingly flat.



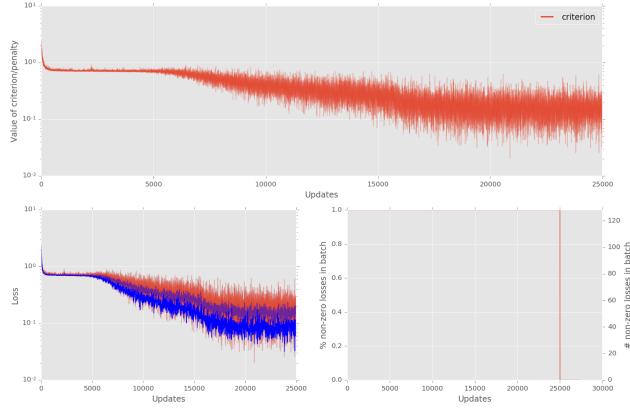(a) Training-log of the loss and active triplet count.
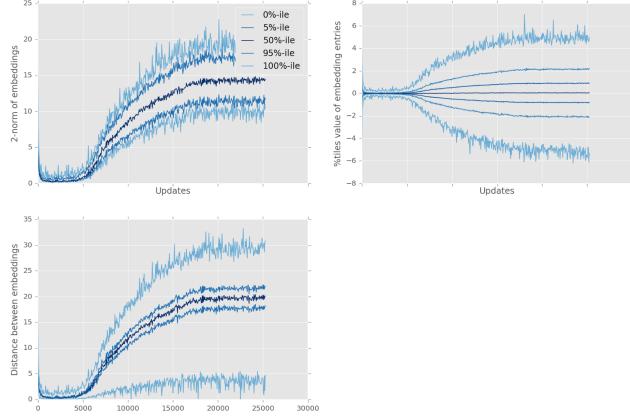


(b) Training-log of the embeddings in the minibatch.

Figure 6: Training-log of LuNet on Market1501 using the *batch hard* triplet loss with soft margin. The embeddings keep moving apart as even the loss shows a steady downward trend.

Second, each Subfigure (b) (blue plots), monitors statistics about the embeddings computed during training. Different lines show 0, 5, 50, 95, and 100-percentiles within a mini-batch, thus visualizing the distribution of values. The top-left plot, "2-norm of embeddings", shows the norms of the embeddings, thus visualizing whether the embedding-space shrinks towards 0 or expands. The top-right plot, "%tiles value of embedding entries" shows these same statistics over the individual numeric entries in the embedding vectors. The only use we found for this plot is noticing when embeddings collapse to all-zeros vs. some other value. Finally, the bottom-left plot, "Distance between embeddings", is the most revealing, as it shows the same percentiles over all pairwise distances between the embeddings within a mini-batch. Due to a bug, the x-axis is unfortunately mislabeled in some cases.

Let us now start by looking at the logs of two very suc-

cessful runs: the LuNet training from scratch on Market-1501 with the *batch hard* loss with margin $0.2$ and in the soft-margin formulation, see Fig. 5 and Fig. 6, respectively. The first observation is that, although they reach similar final scores, they learn significantly different embeddings. Looking at the embedding distances and norms, it is clear that the soft-margin formulation keeps pushing the embeddings apart, whereas the $0.2$-margin keeps the norm of embeddings and their distances bounded. The effect of exponentially decaying the learning-rate is also clearly visible: starting at $15\,000$ updates, both the loss as well as the number of non-zero entries in a mini-batch start to strongly decrease again, before finally converging from $20\,000$ to $25\,000$ updates.

A network getting stuck only happened to us with too weak network architectures, or when using offline hard mining (OHM), the latter can be seen in Fig 8.

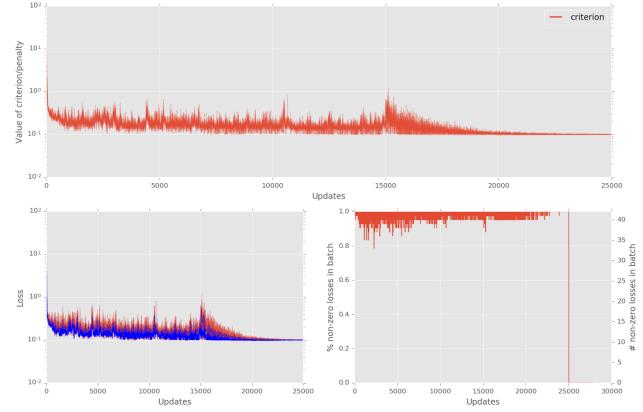(a) Training-log of the loss and active triplet count.



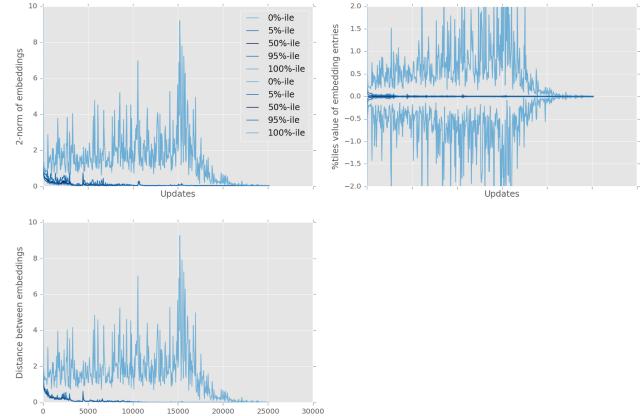(a) Training-log of the loss and active triplet count.



(b) Training-log of the embeddings in the minibatch.

Figure 7: Training-log of a very small network on Market-1501 using the *batch hard* triplet loss with soft margin. The difficult "packed" phase is clearly visible.



(b) Training-log of the embeddings in the minibatch.

Figure 8: Training-log of LuNet on MARS when using of-fline hard mining (OHM) with margin $0.1$. This is one of the runs that collapsed and never got past the difficult phase.

## E. Extended Comparison Tables

We show extended versions of the two state-of-the-art comparison tables in the main paper. We add additional methods that were left out due to space reasons, or because the approaches are orthogonal to ours. The latter could be integrated with our approach in a straightforward manner. Market-1501 and Mars results are shown in Table 7 and CUHK03 results in Table 8.

Next, let us turn to Fig. 7, which shows the training-logs of a very small net (not further specified in this paper). We can clearly see that the network first pulls all embed-dings towards their center of gravity, as evidenced by the quickly decreasing embedding norms, entries, as well as distances in the first few hundred updates. (More visible when zooming-in on a computer.) Once they are all close to each other, the networks really struggles to make them all "pass each other" to reach "their" clusters. As soon as this difficult phase is overcome, the embeddings are spread around the space to quickly decrease the loss. This is where the training becomes "fragile" and prone to collapsing: if the embeddings never pass each other, training gets stuck. This behavior can also be observed in Figures 5 and 6, al-though to a much lesser extent, as the network is powerful enough to quickly overcome this difficult phase.

| | Type | Market-1501 SQ | | | Market-1501 MQ | | | MARS | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | mAP | rank-1 | rank-5 | mAP | rank-1 | rank-5 | mAP | rank-1 | rank-5 |
| TriNet | E | 69.14 | 84.92 | **94.21** | 76.42 | 90.53 | **96.29** | **67.70** | **79.80** | **91.36** |
| LuNet | E | 60.71 | 81.38 | 92.34 | 69.07 | 87.11 | 95.16 | 60.48 | 75.56 | 89.70 |
| IDE (R) + ML ours | I | 58.06 | 78.50 | 91.18 | 67.48 | 85.45 | 94.12 | 57.42 | 72.42 | 86.21 |
| LOMO + Null Space [45] | E | 29.87 | 55.43 | - | 46.03 | 71.56 | - | - | - | - |
| Gated siamese CNN [39] | V | 39.55 | 65.88 | - | 48.45 | 76.04 | - | - | - | - |
| Spindle Net [48] | I | - | 76.9 | 91.5 | - | - | - | - | - | - |
| CAN [25] | E | 35.9 | 60.3 | - | 47.9 | 72.1 | - | - | - | - |
| SSM [3] | - | 68.80 | 82.21 | - | 76.18 | 88.18 | - | - | - | - |
| JLML [22] | I | 65.5 | 85.1 | - | 74.5 | 89.7 | - | - | - | - |
| SVDNet [35] | I | 62.1 | 82.3 | - | - | - | - | - | - | - |
| CNN + DCGAN [54] | I | 56.23 | 78.06 | - | 68.52 | 85.12 | - | - | - | - |
| DPFL [7] | I | **73.1** | **88.9** | - | **80.7** | **92.3** | - | - | - | - |
| ResNet 50 (I+V)† [52] | I+V | 59.87 | 79.51 | 90.91 | 70.33 | 85.84 | 94.54 | - | - | - |
| MobileNet+DML† [47] | I | 68.83 | 87.73 | - | 77.14 | 91.66 | - | - | - | - |
| DTL† [10] | E | 41.5 | 63.3 | - | 49.7 | 72.4 | - | - | - | - |
| DTL† [10] | I+V | 65.5 | 83.7 | - | 73.8 | 89.6 | - | - | - | - |
| APR (R, 751)† [24] | I | 64.67 | 84.29 | 93.20 | - | - | - | - | - | - |
| PAN† [53] | I | 63.35 | 82.81 | - | 71.72 | 88.18 | - | - | - | - |
| IDE (C) + ML [49] | I | - | - | - | - | - | - | 47.6 | 65.3 | 82.0 |
| Latent Parts (Fusion) [20] | I | 57.53 | 80.31 | - | 66.70 | 86.79 | - | 56.05 | 71.77 | 86.57 |
| IDE (R) + ML [55] | I | 49.05 | 73.60 | - | - | - | - | 55.12 | 70.51 | - |
| Spatial-Temporal RNN [57] | E | - | - | - | - | - | - | 50.7 | 70.6 | 90.0 |
| CNN + Video† [46] | I | - | - | - | - | - | - | - | 55.5 | 70.2 |
| TriNet (Re-ranked) | E | **81.07** | **86.67** | 93.38 | **87.18** | **91.75** | 95.78 | **77.43** | **81.21** | **90.76** |
| LuNet (Re-ranked) | E | 75.62 | 84.59 | 91.89 | 82.61 | 89.31 | 94.48 | 73.68 | 78.48 | 88.74 |
| IDE (R) + ML ours (Re-ra.) | I | 71.38 | 81.62 | 89.88 | 79.78 | 86.79 | 92.96 | 69.50 | 74.39 | 85.86 |
| IDE (R) + ML (Re-ra.) [55] | I | 63.63 | 77.11 | - | - | - | - | 68.45 | 73.94 | - |
| PAN (Re-ra.)† [53] | I | 76.65 | 85.78 | - | 83.79 | 89.79 | - | - | - | - |

Table 7: Scores on both the Market-1501 and MARS datasets. The top and middle contain our scores and those of the current state-of-the-art respectively. The bottom contains several methods with re-ranking [55]. The different types represent the optimization criteria, where I stands for identification, V for verification and E for embedding. All our scores include test-time augmentation. The best scores are bold. †: Concurrent work only published on arXiv.

| | Type | Labeled | | Detected | |
|---|---|---|---|---|---|
| | | r-1 | r-5 | r-1 | r-5 |
| TriNet | E | **89.63** | **99.01** | **87.58** | **98.17** |
| Gated siamese CNN [39] | V | - | - | 61.8 | 86.7 |
| DGD* [42] | I | 75.3 | - | - | - |
| LOMO + Null Space [45] | E | 62.55 | 90.05 | 54.70 | 84.75 |
| SSM [3] | - | 76.6 | 94.6 | 72.7 | 92.4 |
| CAN [25] | E | 77.6 | 95.2 | 69.2 | 88.5 |
| Latent Parts (Fusion) [20] | I | 74.21 | 94.33 | 67.99 | 91.04 |
| Spindle Net* [48] | I | 88.5 | 97.8 | - | - |
| JLML [22] | I | 83.2 | 98.0 | 80.6 | 96.9 |
| SVDNet [35] | I | - | - | 81.8 | - |
| CNN + DCGAN [54] | I | - | - | 84.6 | 97.6 |
| DPFL [7] | I | 82.8 | - | 82.0 | - |
| DTL$^{\dagger}$ [10] | I+V | 85.4 | - | 84.1 | - |
| ResNet 50 (I+V)$^{\dagger}$ [52] | I+V | - | - | 83.4 | 97.1 |

Table 8: Scores on CUHK03 for TriNet and a set of recent top performing methods. The best scores are highlighted in bold. $^{\dagger}$: Concurrent work only published on arXiv. *: The method was trained on several additional datasets.

## F. LuNet's Architecture

The details of the LuNet architecture for training from scratch can be seen in Table 9. The input image has three channels and spatial dimensions $128\times64$. Most *Res-blocks* are of the "bottleneck" type [15], meaning for given numbers $n_1, n_2, n_3$ in the table, they consist of a $1\times1$ convolution from the number of input channels $n_1$ to the number of intermediate channels $n_2$, followed by a $3\times3$ convolution keeping the number of channel constant, and finally another $1\times1$ convolution going from $n_2$ channels to $n_3$ channels. Only the last Res-block, whose exact filter sizes are given in the table, is an exception to this. All ReLUs, including those in Res-blocks, are leaky [27] by a factor of $0.3$; although we do not have side-by-side experiments comparing the benefits, we expect them to be minor. All convolutional weights are initialized following He *et al.* [13], whereas we initialized the final *Linear* layers following Glorot *et al.* [11]. Batch-normalization [16] is essential to train such a network, and makes the exact initialization less important.

| Type | Size |
|---|---|
| Conv | $128\times7\times7\times3$ |
| Res-block | $128, 32, 128$ |
| MaxPool | pool $3\times3$, stride $(2\times2)$, padding $(1\times1)$ |
| Res-block | $128, 32, 128$ |
| Res-block | $128, 32, 128$ |
| Res-block | $128, 64, 256$ |
| MaxPool | pool $3\times3$, stride $(2\times2)$, padding $(1\times1)$ |
| Res-block | $256, 64, 256$ |
| Res-block | $256, 64, 256$ |
| MaxPool | pool $3\times3$, stride $(2\times2)$, padding $(1\times1)$ |
| Res-block | $256, 64, 256$ |
| Res-block | $256, 64, 256$ |
| Res-block | $256, 128, 512$ |
| MaxPool | pool $3\times3$, stride $(2\times2)$, padding $(1\times1)$ |
| Res-block | $512, 128, 512$ |
| Res-block | $512, 128, 512$ |
| MaxPool | pool $3\times3$, stride $(2\times2)$, padding $(1\times1)$ |
| Res-block | $512\times(3\times3\times512), 128\times(3\times3\times512)$ |
| Linear | $1024\times512$ |
| Batch-Norm | $512$ |
| ReLU | |
| Linear | $512\times128$ |

Table 9: The architecture of LuNet.

Figure 9: Barnes-Hut t-SNE [38] of our learned embeddings for the Market-1501 test-set. Best viewed when zoomed-in.

## G. Full t-SNE Visualization

Figure 9 shows the full Barnes-Hut t-SNE visualization from which the teaser image (Fig. 1 in the paper) was cropped. We used a subset of 6000 images from the Market-1501 test-set and a perplexity of 5000 for this visualization.