

Human-In-The-Loop Person Re-Identification

Hanxiao Wang, Shaogang Gong, Xiatian Zhu, Tao Xiang

School of EECS, Queen Mary University of London, UK
{hanxiao.wang, s.gong, xiatian.zhu, t.xiang}@qmul.ac.uk

Abstract. Current person re-identification (re-id) methods assume that (1) pre-labelled training data are available for every camera pair, (2) the gallery size for re-identification is moderate. Both assumptions scale poorly to real-world applications when camera network size increases and gallery size becomes large. Human verification of automatic model ranked re-id results becomes inevitable. In this work, a novel human-in-the-loop re-id model based on Human Verification Incremental Learning (HVIL) is formulated which does not require any pre-labelled training data to learn a model, therefore readily scalable to new camera pairs. This HVIL model learns cumulatively from human feedback to provide instant improvement to re-id ranking of each probe on-the-fly enabling the model scalable to large gallery sizes. We further formulate a Regularised Metric Ensemble Learning (RMEL) model to combine a series of incrementally learned HVIL models into a single ensemble model to be used when human feedback becomes unavailable.

Keywords: Person re-identification; incremental learning; human-in-the-loop; metric ensemble.

1 Introduction

State-of-the-art person re-identification (re-id) models are dominated by supervised learning approaches [1–12], which employ a *train-once-and-deploy* scheme (Fig. 1(a)). That is, a pre-labelled training data set with given cross-view true-matching identities is first collected and used to learn a model. The learned model is then deployed to new data without any modification. Based on this approach, the re-id community has witnessed over the past two years ever-increased re-id matching accuracy on increasingly larger sized benchmarks with more identities. For instance, the CUHK03 benchmark [9] contains 13,164 images of 1,360 identities which is significantly larger than the early VIPeR [13] and iLIDS [14] benchmarks. The state-of-the-art Rank-1 matching accuracy on CUHK03 is now in 50-60% [12], doubling the best performance reported merely a year ago [9].

One inevitable question arises: Are we close to an automated re-id solution capable of deployment in the real-world? The answer is no. This is because existing supervised learning based re-id methods make two critical assumptions, both of which are invalid in the real-world (unscalable): (1) A manually pre-labelled pairwise training data set is assumed available for every camera pair. However, this is neither scalable (prohibitive to collect in the real-world as there are

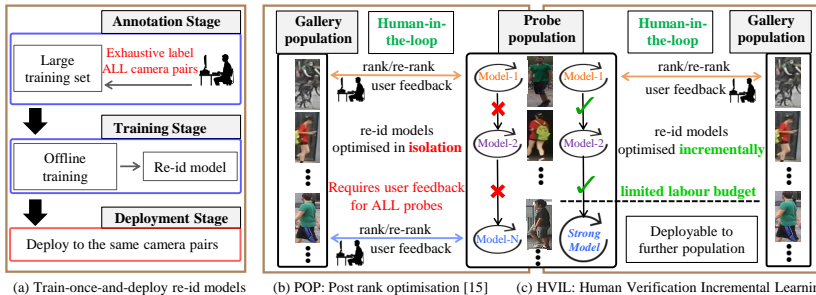


Fig. 1. (a) Conventional *train-once-and-deploy* re-id strategy requires pre-labelled training data collection. (b) POP [15]: A recent *human-in-the-loop re-id* approach which optimises probe-specific models in isolation. (c) HVIL: The proposed new incremental human-in-the-loop re-id model.

quadratic number of camera pairs), nor plausible (there may not exist sufficiently large number of training people reappearing in every pair of camera views). (2) The size of the training dataset is assumed either significantly greater or no less than that of test gallery population on which the learned model will be deployed. For instance, given the standard splits of the CUHK03 benchmark, the training set consists of paired images of 1,260 person identities from six different camera views (on average 4.8 image samples per person per camera view), whilst the test gallery set consists of only 100 identities each with a single image (one-shot setting). The test set’s identity size is thus 10 times less than that of the training set, and has approximately 50 times less images. In a real-world, the size of any deployment gallery population is almost always much greater than any pre-labelled training data size even if such training data were made available. In a public space such as an underground station, there are easily over 1,000 people passing through a camera network every hour, resulting in a typical gallery population size of over 10,000 in a day. It was observed from our experiments that a 10-fold increase in gallery size leads to a 10-fold decrease in re-id Rank-1 performance, resulting in a single-digit Rank-1 score, even when the state-of-the-art re-id models were trained from sufficiently sized labelled data. Given such single-digit Rank-1 scores, human operators are required to verify any true match given a probe from a rather large rank list.

To overcome the inherent limitations of the two aforementioned assumptions from pre-labelling based supervised learning, an attractive alternative approach is to explore *human-in-the-loop* for person re-identification (Fig. 1(b)). Such an approach is inherently more scalable compared to conventional pre-trained re-id models because it does not assume the collection of pre-labelled training data. Human-in-the-loop verification can be considered as a form of “labelling effort”. However, this on-the-fly verification approach has two significant advantages over the conventional approach that requires pre-labelling data for training: (1) It requires much less labelling-effort (the number of feedback data from human verification is typically in tens rather than thousands required for pre-labelling training data); (2) It focuses on optimising the re-id ranking of each probe directly in the test gallery population, rather than learning a distance metric in a separate training set and blindly assuming its adaptability to the test gallery population.

In this work we develop a re-id model without the need for pre-labelled training. Crucially, it can be improved incrementally by human verification and benefits from more flexible human feedback (similar/dissimilar). As a result, it enables a human to re-id rapidly a given probe image after only a handful of feedback verifications even when the gallery size is large. More specifically, a Human Verification Incremental Learning (HVIL) model (Fig. 1(c)) is formulated to maximise the effectiveness of human-in-the-loop feedback by incorporating: (1) *Flexible feedback* - HVIL allows for weak human feedback (similar/dissimilar) without the need for exhaustive user search in the ranked list, instead of being restricted to only true/false verifications. (2) *Immediate benefit* - By introducing a new online incremental distance metric learning model, HVIL enables real-time response to human feedback by rapidly presenting a freshly optimised ranking list. (3) *The older the wiser* - HVIL is updated cumulatively on-the-fly utilising multiple user feedback per probe and optimised incrementally for each new probe given what been learned from all previous probes. (4) *A strong ensemble model* - An additional Regularised Metric Ensemble Learning (RMEL) model is introduced by taking all the incrementally optimised per-probe models as a set of “weak” models [16, 17] and constructing a “strong” ensemble model for performing re-id tasks when human feedback becomes unavailable.

Related Work - Current best performing person re-id methods are fully supervised but they require a large number of pre-labelled training data from every camera pair for building camera-pair specific distance metric models [18, 2–6, 10, 7–9, 19, 20, 11, 12]. Their usefulness and scalability are inherently limited in real-world applications especially with large camera networks. This problem becomes more acute for the more recent deep learning based methods [21, 22, 19, 23, 9] which need more labelled training data to function. To relax this need for labelling, existing attempts include semi-supervised [24, 25], unsupervised [26–28], and transfer learning [29–32]. However, all of these strategies are weak in performance compared to fully supervised learning - without labelled data, they are unable to learn strong discriminative information for cross-view people re-identification. In contrast, the proposed HVIL model learns discriminatively from human feedback instead of pre-labelled image pairs, and is capable of yielding much superior person re-id matching accuracy than the state-of-the-art supervised re-id models, with added advantages of costing much less human feedback as “labelling effort” and being more scalable to large test gallery sizes.

Very few human-in-the-loop re-id methods were reported before, nor received much attention. Abir et al. [33] assumed a pre-labelled training set available per person *in addition* to human-in-the-loop verification. Hirzer et al. [34] considered a form of human feedback which is ill-posed: It only allows a user to verify whether a *true* match is within the top-N ranking list. This limits significantly the effectiveness of human feedback and can waste expensive human labour when a true match cannot be found in the top-N ranks. More recently, Liu et al. [15] proposed the POP model (Fig. 1(b)), which allows a user to identify correct matches more rapidly and accurately by accommodating more flexible feedback information. However, both [34, 15] are limited inherently due to the fact that

they treat each probe as an independent retrieval task, i.e. the process of learning a model for each probe does not benefit learning models for other probes. This lack of improving model-learning cumulatively with increasing human feedback is both suboptimal and in danger of disengaging the human in the loop. In contrast, the proposed HVIL re-id framework (Fig. 1(c)) enables incremental model improvement from cumulative human feedback. Moreover, the proposed RMEL ensemble model further benefits from previous human verification effort even when human feedback is no long available.

Contributions - (1) We formulate a new approach to person re-id for a model to be optimised cumulatively by human feedback on-the-fly with each re-id task at hand without pre-labelled training and being effective for large gallery sizes. (2) A Human Verification Incremental Learning (HVIL) model is introduced for distance metric optimisation by flexible human feedback continuously in real-time and from cumulative feedback when more probe images are searched. (3) A Regularised Metric Ensemble Learning (RMEL) model is constructed for a strong ensemble model when human feedback becomes unavailable. The advantages of the proposed approach is validated by extensive comparisons against contemporary image retrieval methods and state-of-the-art supervised person re-id models on two largest re-id benchmarks CUHK03 [9] and Market-1501 [35].

2 Human-in-the-Loop Incremental Learning

2.1 Problem Formulation

Suppose an image is denoted by a feature vector $\mathbf{x} \in \mathbb{R}^d$. The *human-in-the-loop re-id* problem is formulated as: (1) For each image \mathbf{x}^p in a probe set $\mathcal{P} = \{\mathbf{x}_i^p\}_{i=1}^{N_p}$, \mathbf{x}^p is matched against a gallery set $\mathcal{G} = \{\mathbf{x}_i^g\}_{i=1}^{N_g}$ and an initial ranking list is generated by a re-id ranking function $f(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$, according to ranking scores $f_{\mathbf{x}^p}(\mathbf{x}_i^g)$. (2) A human operator (user) browses the gallery ranking list to verify the existence and the rank of any true match for \mathbf{x}^p . Human feedback is generated when a ranked gallery image \mathbf{x}^g is selected by the user with a label $y \in \{\text{true, dissimilar, similar}\}$. Once a feedback on probe \mathbf{x}^p is received, parameters of $f(\cdot)$ are updated instantly to re-order the gallery ranking list and give the user immediate reward for the feedback. (3) When either a true match is found or a pre-determined maximum round of feedback is reached, the next probe is presented for re-id in the gallery set. In contrast to pre-labelling training data required by conventional *train-once-and-deploy* re-id schemes, *human-in-the-loop re-id* has two unique characteristics: (a) Due to human patience and limited labour budget [34], a user is only interested in the top ranked gallery images, and a user’s feedback on each probe is limited. (b) Rather than verifying *only* true (positive) matches in the gallery for each probe, which are inherently very few if any among the top ranks¹, it is a much easier and more rewarding task for the user to give feedback on the many top ranked negative gallery instances: *strong-negative* (dissimilar) - “*definitely not the one I am looking for*”, and *weak-negative* (similar) - “*looks similar but not the same person*” [15].

¹ In a large size gallery set, true matches are often scarce (only one-shot) and overwhelmed (appear in low-ranks) by false matches of high-ranks in the rank list.

2.2 Modelling Human Feedback as a Loss Function

Formally, we wish to construct an incrementally optimised ranking function, $f_{\mathbf{x}^p}(\mathbf{x}_i^g) : \mathbb{R}^d \rightarrow \mathbb{R}$, where $f(\cdot)$ can be estimated by three types of human feedback $y \in L = \{m, s, w\}$ as *true-match*, *strong-negative*, and *weak-negative* respectively. Inspired by [36–38], we define a ranking error ($\mathcal{L}_{\text{loss}}$) function for a feedback y on a human selected gallery sample \mathbf{x}^g given a probe \mathbf{x}^p as:

$$\text{err}(f_{\mathbf{x}^p}(\mathbf{x}^g), y) = \mathcal{L}_y(\text{rank}(f_{\mathbf{x}^p}(\mathbf{x}^g))), \quad (1)$$

where $\text{rank}(f_{\mathbf{x}^p}(\mathbf{x}^g))$ denotes the rank of \mathbf{x}^g given by $f_{\mathbf{x}^p}(\cdot)$, defined as:

$$\text{rank}(f_{\mathbf{x}^p}(\mathbf{x}^g)) = \sum_{\mathbf{x}_i^g \in \mathcal{G} \setminus \mathbf{x}^g} \mathcal{I}(f_{\mathbf{x}^p}(\mathbf{x}_i^g) \geq f_{\mathbf{x}^p}(\mathbf{x}^g)), \quad (2)$$

where $\mathcal{I}(\cdot)$ is the indicator function. The loss function $\mathcal{L}_y(\cdot) : \mathbb{Z}^+ \rightarrow \mathbb{R}^+$ transforms a rank into a loss. We introduce a novel re-id ranking loss defined as:

$$\mathcal{L}_y(k) = \begin{cases} \sum_{i=1}^k \alpha_i, & \text{if } y \in \{m, w\} \\ \sum_{i=k+1}^{n_g} \alpha_i, & \text{if } y \in \{s\} \end{cases}, \quad \text{with } \alpha_1 \geq \alpha_2 \geq \dots \geq 0. \quad (3)$$

Note, different choices of α_i lead to specific model responses to human feedback. We set $\alpha_i = \frac{1}{i}$ (large penalty with steep slope) when y indicates a *true-match*, and $\alpha_i = \frac{1}{n_g - i + 1}$ with n_g the gallery size (small penalty with gentle slope) when y represents a *weak-negative* or *strong-negative*. Such a ranking loss is designed to favour a model update behaviour so that: (1) *true-matches* are quickly pushed up to the top ranks, whilst (2) *weak-/strong-negatives* are mildly moved towards the top/bottom rank direction. Our experiments (Sec. 4.1) show that such a ranking loss criterion boosts very effectively the Rank-1 matching rate and pushes quickly *true-matches* to the top ranks at each iteration of human feedback.

2.3 Real-time Model Update for Instant Feedback Reward

Given the re-id ranking loss function defined in Eqn. (3), we wish to have real-time model update to human feedback therefore providing instant reward to user labour effort. To that end, we consider the re-id ranking model $f(\cdot)$ as a negative Mahalanobis distance metric:

$$f_{\mathbf{x}^p}(\mathbf{x}^g) = - [(\mathbf{x}^p - \mathbf{x}^g)^\top \mathbf{M}(\mathbf{x}^p - \mathbf{x}^g)], \quad \mathbf{M} \in S_+^d. \quad (4)$$

The positive semi-definite matrix \mathbf{M} consists of model parameters to be learned.

Knowledge cumulation by online learning - In previous works [34, 15], $f(\cdot)$ is only optimised in isolation for each probe without benefiting from previous feedback on other probes. To overcome this limitation, we wish to optimise $f(\cdot)$ incrementally in an online manner [39] for maximising the value of limited human feedback labour budget. Moreover, to achieve real-time human-in-the-loop feedback and reward, $f(\cdot)$ needs to be estimated on each human feedback.

Formally, given a new probe \mathbf{x}_t^p at time step $t \in \{1, \dots, \tau\}$ (τ the pre-defined budget), a user is presented with a gallery rank list computed by the previously estimated model \mathbf{M}_{t-1} instead of re-initialising a new ranking function from scratch for this new probe. The user then verifies a gallery image \mathbf{x}_t^g in the top ranks with a label y_t , generating a labelled triplet $(\mathbf{x}_t^p, \mathbf{x}_t^g, y_t)$. Given Eqn. (3), this triplet has a corresponding loss as $\mathcal{L}^{(t)} = \mathcal{L}_{y_t}(\text{rank}(f_{\mathbf{x}_t^p}(\mathbf{x}_t^g)))$. We update the ranking model by minimising the following object function:

$$\mathbf{M}_t = \underset{\mathbf{M} \in S_{\perp}^d}{\text{argmin}} \Delta_F(\mathbf{M}, \mathbf{M}_{t-1}) + \eta \mathcal{L}^{(t)}, \quad (5)$$

where Δ_F is a Bregman divergence measure, defined by an arbitrary differentiable convex function F , for regularising the discrepancy between \mathbf{M} and \mathbf{M}_{t-1} . The set S_{\perp}^d defines a PSD cone, and the tradeoff parameter $\eta > 0$ balances the model update divergence and empirical loss. This optimisation updates incrementally the ranking model adopted from the previous probe by encoding user feedback on the current probe.

Loss approximation for real-time optimisation - In order to encourage and maintain user engagement in verification feedback, real-time online incremental metric learning is required. However, as $\mathcal{L}^{(t)}$ is discontinuous, the overall objective function cannot be optimised efficiently by gradient-based methods. We thus approximate the loss function by a continuous upper bound [36] so that it is differentiable w.r.t. \mathbf{M} :

$$\tilde{\mathcal{L}}^{(t)} = \frac{1}{\mathcal{N}_t^-} \sum_{\mathbf{x}_i^g \in \mathcal{G} \setminus \mathbf{x}_t^g} \mathcal{L}_{y_t} \left(\text{rank} \left(f_{\mathbf{x}_t^p}(\mathbf{x}_i^g | \mathbf{M}_{t-1}) \right) \right) h_{y_t} \left(f_{\mathbf{x}_t^p}(\mathbf{x}_i^g | \mathbf{M}_t) - f_{\mathbf{x}_t^p}(\mathbf{x}_t^g | \mathbf{M}_t) \right)^2, \quad (6)$$

where $f_{\mathbf{x}_t^p}(\mathbf{x}_i^g | \mathbf{M}_{t-1})$ denotes the function value of $f_{\mathbf{x}_t^p}(\mathbf{x}_i^g)$ parametrised by \mathbf{M}_{t-1} , and $h_{y_t}(\cdot)$ represents a hinge loss function defined as:

$$h_{y_t}(f_{\mathbf{x}_t^p}(\mathbf{x}_i^g) - f_{\mathbf{x}_t^p}(\mathbf{x}_t^g)) = \begin{cases} \max(0, 1 - f_{\mathbf{x}_t^p}(\mathbf{x}_i^g) + f_{\mathbf{x}_t^p}(\mathbf{x}_t^g)), & \text{if } y_t \in \{m, w\} \\ \max(0, 1 - f_{\mathbf{x}_t^p}(\mathbf{x}_i^g) + f_{\mathbf{x}_t^p}(\mathbf{x}_t^g)), & \text{if } y_t \in \{s\} \end{cases}. \quad (7)$$

The normaliser \mathcal{N}_t^- in Eqn. (6) is the amount of violators, i.e. the gallery instances that generate non-zero hinge loss in Eqn. (7) w.r.t. triplet $(\mathbf{x}_t^p, \mathbf{x}_t^g, y_t)$.

Learning speed-up by most violator update - Given the approximation in Eqn. (6), we can exploit the stochastic gradient descent (SGD) algorithm [40] for optimising Eqn. (5) by iteratively updating on sub-sampled batches of all violators. However, the computational overhead of iterative updates can be large due to possibly many violators, and thus not meeting the real-time requirement. To address this problem, we explore a *most violator update* strategy, that is, to perform metric updates using *only* the violator \mathbf{x}_v^g with the most violation (Eqn. (7)). The final approximated empirical loss is then estimated as:

$$\tilde{\mathcal{L}}_v^{(t)} = \mathcal{L}_{y_t} \left(\text{rank} \left(f_{\mathbf{x}_t^p}(\mathbf{x}_i^g | \mathbf{M}_{t-1}) \right) \right) h_{y_t} \left(f_{\mathbf{x}_t^p}(\mathbf{x}_i^g | \mathbf{M}_t) - f_{\mathbf{x}_t^p}(\mathbf{x}_v^g | \mathbf{M}_t) \right)^2. \quad (8)$$

By replacing $\mathcal{L}^{(t)}$ in Eqn. (5) with $\tilde{\mathcal{L}}_v^{(t)}$, and setting the gradient of Eqn. (5) to zero, we yield the following ranking metric online update criterion:

$$\mathbf{M}_t = g^{-1} \left(g(\mathbf{M}_{t-1}) - \eta \nabla_{\mathbf{M}} \tilde{\mathcal{L}}_v^{(t)} \right), \quad (9)$$

where $g(\cdot)$ denotes the derivative of F (Eqn. (5)) w.r.t. \mathbf{M} [41]. For the form of F , we adopt Burg matrix divergence [42]:

$$\Delta_F(\mathbf{M}, \mathbf{M}_{t-1}) = \text{tr}(\mathbf{M}\mathbf{M}_{t-1}^{-1}) - \log \det(\mathbf{M}\mathbf{M}_{t-1}^{-1}). \quad (10)$$

Eqn. (9) can be readily optimised by any gradient-based update schemes [43, 41]. We adopted the LogDet Exact Gradient Online (LEGO) algorithm [44]. This is desirable because Eqn. (9) is solved with a computational complexity of $\mathcal{O}(d^2)$ where d is the feature vector dimension. This avoids eigenvector computation with a cost of $\mathcal{O}(d^3)$ required by most other schemes. Given all the components described above, our final model for Human Verification Incremental Learning (HVIL) enables real-time incremental person re-id model learning with human-in-the-loop feedback. Our extensive experiments (Sec. 4.1) show that this HVIL model provides the fastest human-in-the-loop feedback-reward cycle over other competitors. An overview of the HVIL model is given in Algorithm 1.

Algorithm 1: Human Verification Incremental Learning (HVIL)

Data: Unlabelled probe set \mathcal{P} and gallery set \mathcal{G} ;
Result: Per probe optimised ranking lists; re-id models $\{\mathbf{M}_t\}_{t=1}^{\tau}$;
 Initialisation: $\mathbf{M}_0 = \mathbf{I}$ (identity matrix, equivalent to L_2 distance)
while $t < \tau$ **do**
 Present the next probe $\mathbf{x}_t^p \in \mathcal{P}$;
 for $iter = 1 : \text{maxIter}$ **do**
 // **maxIter:** maximum interaction rounds per probe
 Rank \mathcal{G} with \mathbf{M}_{t-1} against \mathbf{x}_t^p (Eqn. (4));
 Collect human feedback (\mathbf{x}_t^g, y_t) ;
 Locate the most violator \mathbf{x}_v^g and calculate $\tilde{\mathcal{L}}_v^{(t)}$ (Eqn. (7) and Eqn. (8));
 $\mathbf{M}_t = \text{update}(\mathbf{M}_{t-1}, \tilde{\mathcal{L}}_v^{(t)})$ (Eqn. (9)), $t = t + 1$;
 end
end
 Return $\{\mathbf{M}_t\}_{t=1}^{\tau}$.

3 Metric Ensemble Learning for Automated Re-id

Finally, we consider a situation when limited human labour budget is exhausted at time τ and an automated re-id strategy is required for any further probes. In this case, as the HVIL re-id model is optimised incrementally, the model \mathbf{M}_τ optimised by the human verified probe at time τ can be directly deployed. However, it is desirable to construct an even “stronger” model based on metric ensemble learning. Specifically, a side-product of HVIL is a series of models incrementally optimised *locally* for a set of probes with human feedback. We consider them as a set of *globally* “weak” models $\{\mathbf{M}_j\}_{j=1}^{\tau}$, and wish to construct a *single globally strong model* for re-id further probes without human feedback.

Regularised Metric Ensemble Learning - Given weak models $\{\mathbf{M}_j\}_{j=1}^\tau$, we compute a distance vector $\mathbf{d}_{ij} \in \mathbb{R}^\tau$ for any probe-gallery pair $(\mathbf{x}_j^g, \mathbf{x}_i^p)$:

$$\mathbf{d}_{ij} = - \left[f_{\mathbf{x}_i^p}(\mathbf{x}_j^g | \mathbf{M}_1), \dots, f_{\mathbf{x}_i^p}(\mathbf{x}_j^g | \mathbf{M}_t), \dots, f_{\mathbf{x}_i^p}(\mathbf{x}_j^g | \mathbf{M}_\tau) \right]^\top. \quad (11)$$

The objective of metric ensemble learning is to obtain an optimal combination of these distances for producing a single globally optimal distance. Here we consider the ensemble ranking function $f_{\mathbf{x}_i^p}^{ens}(\mathbf{x}_j^g)$ in a bi-linear form (shortened as f_{ij}^{ens}):

$$f_{ij}^{ens} = f_{\mathbf{x}_i^p}^{ens}(\mathbf{x}_j^g) = -\mathbf{d}_{ij}^\top \mathbf{W} \mathbf{d}_{ij}, \quad \text{s.t. } \mathbf{W} \in S_+^\tau, \quad (12)$$

with \mathbf{W} being the model parameters capturing the correlations among all the weak model metrics. In this context, previous work such as [20] is a special case of our model when \mathbf{W} is restricted to be diagonal only.

Objective function - To estimate an optimal ensemble weights \mathbf{W} with most identity-discriminative power, we re-use the true matching pairs verified during the human verification procedure (Sec. 2) as ‘‘training data’’: $\mathcal{X}_{tr} = \{(\mathbf{x}_i^p, \mathbf{x}_i^g)\}_{i=1}^{N_t}$, and their corresponding person identities are denoted by $\mathcal{C} = \{c_i\}_{i=1}^{N_t}$. Note, ‘‘training data’’ here are only for estimating ensemble weights, not for learning a distance metric. Since the ranking score f_{ij}^{ens} in Eqn. (12) is either negative or zero, we consider that in the extreme case, an *ideal* ensemble function f^* should provide the following ranking scores : $f_{ij}^* = 0$ for $c_i = c_j$, and $f_{ij}^* = -1$ for $c_i \neq c_j$. Using \mathbf{F}^* to denote such an ideal ranking score matrix and \mathbf{F}^{ens} to denote an estimated score matrix by a given \mathbf{W} with Eqn. (12), our proposed objective function for metric ensemble learning is then defined as:

$$\min_{\mathbf{W}} \|\mathbf{F}^{ens} - \mathbf{F}^*\|_F^2 + \nu \mathcal{R}(\mathbf{W}), \quad \text{s.t. } \mathbf{W} \in S_+^\tau, \quad (13)$$

where $\|\cdot\|_F$ denotes a Frobenius norm, and $\mathcal{R}(\mathbf{W})$ a regulariser on \mathbf{W} with parameter ν controlling the regularisation strength. Whilst common choices of $\mathcal{R}(\mathbf{W})$ include L_1 , Frobenius norm, or matrix trace, we introduce the following regularisation for a Regularised Metric Ensemble Learning (RMEL) re-id model:

$$\mathcal{R}(\mathbf{W}) = - \sum_{i,j} f_{ij}^{ens}, \quad \text{for } c_i = c_j. \quad (14)$$

Our intuition is to impose severe penalties for true match pairs with low ranking scores since they deliver the most informative discriminative information for cross-view person re-id, whilst false match pairs are either less informative (strong-negative) or non-discriminative (weak-negative).

Optimisation - Eqn. (13) is strictly convex with a guaranteed global optimal so it can be optimised by any off-the-shelf toolboxes [45]. We adopt the standard first-order projected gradient descent algorithm [46]. Given the estimated optimal ensemble weight matrix \mathbf{W} and the weak models $\{\mathbf{M}_j\}_{j=1}^\tau$, a single strong ensemble model (Eqn. (12)) is made available for performing automated re-id of any further probes on the gallery population. Our experiments (Sec. 4.2) show that the proposed RMEL algorithm achieves superior performance compared to state-of-the-arts supervised re-id models given the same amount of labelled data.

4 Experiments

Two experiments were conducted: (1) The proposed HVIL model was evaluated under a *human-in-the-loop re-id* setting and an *enlarged* test gallery population was used to reflect real-world use-cases. (2) In the event of limited human labour budget being exhausted and human feedback becoming unavailable, the proposed HVIL-RMEL model was evaluated under an *automated re-id* setting.

Datasets - Two largest person re-id benchmarks: CUHK03 [9] and Market-1501 [35] were chosen for evaluation due to the need for large test gallery size. CUHK03 contains 13,164 automatically detected bounding boxes of 1,360 people; Market-1501 consists of 32,668 detections of 1,501 people. Both datasets cover 6 outdoor surveillance cameras with severely divergent and unknown viewpoints, illumination conditions, (self)-occlusion and background clutter.

Data partitions - For each dataset, we randomly selected 1,000 identities D_{p1} (p stands for population) as the partition to perform *human-in-the-loop re-id* experiments. The remaining partition of people D_{p2} (360 on CUHK03, and 501 on Market-1501) were separated for evaluating the proposed model against state-of-the-art supervised re-id methods for *automated re-id* (see details in Sec. 4.1 and Sec. 4.2). To obtain statistical reliability, we generated 6 different trials $\{D_{p1}^i, D_{p2}^i\}_{i=1}^6$ for experiments and reported their averaged results.

Visual features - The descriptor of [47] was adopted for person image representation. The feature vector (5,138 dimensions) was a concatenation of colour, HOG [48] and LBP [49] histograms extracted from horizontal rectangular stripes.

4.1 Evaluation on Human-in-the-Loop Person Re-Id

We evaluated the performance of our HVIL model in *human-in-the-loop re-id* setting, along with detailed human feedback statistics analysis.

Human feedback protocol - Human feedback were collected on all 6 trials of D_{p1}^i partitions in 6 independent sessions by 2 volunteers as users, i.e. each trial for one different session. The human labour budget in each session was limited to the maximum of 300 probes. For testing, the standard single-shot re-id scheme [1] is considered, i.e. from the partition D_{p1}^i we selected randomly a single image per identity to form a 300 people/image probe set \mathcal{P}^i and crucially, a much larger 1,000 people/image gallery set \mathcal{G}^i (\mathcal{P}^i and \mathcal{G}^i are from different camera views). During each session, a user was asked to perform *human-in-the-loop re-id* on probes in probe set \mathcal{P}^i against gallery set \mathcal{G}^i . For each probe, a *maximum* of 3 rounds of user interaction are allowed. We limited the users to verify only the top-50 in the rank list (5% of 1,000 gallery set). During each interaction: (1) A user selects one gallery image as either *strong-negative*, *weak-negative*, or *true-match*; and (2) the system takes the feedback, updates the ranking function and returns the re-ordered ranking list, all in real-time (Sec. 2). The HVIL model was evaluated against six existing models for *human-in-the-loop re-id* as follows.

Competitors A - Three existing human-in-the-loop models were compared: (1) POP [15]: The current state-of-the-art *human-in-the-loop re-id* method based on Laplacian SVMs and graph label propagation; (2) Rocchio [50]: A probe

| Dataset | CUHK03 [9] ($N_g = 1000$) | | | | Market-1501 [35] ($N_g = 1000$) | | | |
|--------------|-----------------------------|-------------|-------------|-------------|-----------------------------------|-------------|-------------|-------------|
| | 1 | 50 | 100 | 200 | 1 | 50 | 100 | 200 |
| L2 | 2.9 | 31.1 | 43.2 | 58.2 | 16.1 | 66.6 | 76.6 | 85.0 |
| kLFDA [10] | 5.9 | 47.3 | 60.1 | 75.0 | 21.8 | 85.8 | 91.5 | 96.3 |
| XQDA [11] | 3.7 | 40.2 | 53.6 | 68.5 | 18.3 | 75.1 | 83.5 | 91.1 |
| MLAPG [12] | 4.2 | 39.5 | 52.4 | 66.7 | 24.1 | 84.5 | 91.2 | 95.7 |
| EMR [52] | 46.0 | 47.3 | 51.3 | 60.0 | 53.3 | 64.3 | 75.7 | 85.0 |
| Rocchio [50] | 43.1 | 49.9 | 57.3 | 65.1 | 52.7 | 69.6 | 77.6 | 87.3 |
| POP [15] | 46.3 | 55.7 | 64.0 | 74.3 | 56.0 | 72.7 | 80.6 | 86.3 |
| HVIL (Ours) | 56.1 | 64.7 | 75.7 | 87.4 | 78.0 | 86.0 | 90.3 | 93.4 |

Table 1. Evaluating human-in-the-loop person re-id with CMC performances.

vector modification model updates iteratively the probe’s feature vector based on human feedback, widely used for image retrieval tasks [51]; (3) EMR [52]: A graph-based ranking model that optimises the ranking function by least square regression. For a fair comparison of all four human-in-the-loop models, the users were asked to verify the same probe and gallery data ($\mathcal{P}^i, \mathcal{G}^i$) with three-types of feedback given the ranking-list generated by each model.

Competitors B - In addition, three state-of-the-art conventional supervised person re-id models were also compared: (4) kLFDA [10], (5) XQDA [11], and (6) MLAPG [12]. These supervised re-id methods were trained using fully pre-labelled data in the separate partition D_{p2}^i (CUHK03: averagely 3,483 images of 360 identities; Market-1501: averagely 7,737 images of 501 identities) before being deployed to \mathcal{P}^i (300) and \mathcal{G}^i (1,000) for testing. Note, the underlying human labour effort for pre-labelling the training data to learn these supervised models was significantly greater – exhaustively searching 3,483 and 7,737 *true* matched images respectively for CUHK03 and Market-1501, than that required by the human-in-the-loop methods – between 300 to 900 *indicative* verification (similar, dissimilar, or true) given a maximum of 300 probes on both CUHK03 and Market-1501, so only 1/10th of and weaker user input than supervised models.

Implementation details - For implementing the HVIL model (Sec. 2), the only hyper-parameter η (Eqn. (5)) was set to 0.5 on both datasets. We found that HVIL is insensitive to η with a wide satisfiable range from 10^{-1} to 10^1 . For POP, EMR, and Rocchio, we adopted the authors’ recommended parameter settings as in [15, 50]. For all methods above, we applied L_2 distance as the initial ranking function $f_0(\cdot)$ without loss of generalisation². Note that for HVIL, once $f_0(\cdot)$ was initialised for only the very first probe, it was then optimised incrementally across different probes. In contrast, for POP and EMR and Rocchio, each probe had its own $f_0(\cdot)$ initialised as L_2 since the models are not cumulative across different probes. For supervised methods kLFDA, XQDA and MLAPG, the parameters were determined by cross-validation on D_{p2} with the authors’ published codes. All models adopted the same feature descriptor [47].

Evaluation metrics - Cumulative Match Characteristic (CMC) curves were adopted for performance evaluation. Specifically, we calculated the cumulative recognition rate at each rank position. Expected Rank (ER) is also used for evaluation, defined as the average rank of all true matches. For all human-in-

² No limitation on considering any distance/similarity metrics, either learned or not.

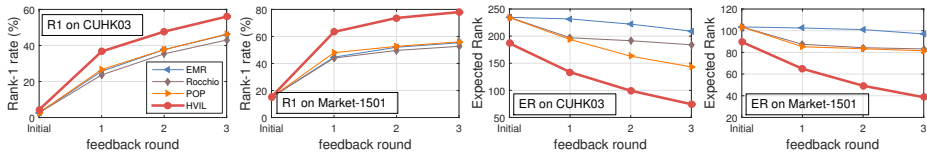


Fig. 2. Comparing Rank-1 score and Expected Rank (ER) on human feedback rounds.

the-loop models, we used the ranking result after the final interaction on each probe for CMC evaluation. The averaged results over all 6 trials are reported.

Comparative results - The person re-id performance of all methods on \mathcal{P}^i and \mathcal{G}^i is shown in Table 1. First, it is evident that when the testing gallery size was enlarged from their standard settings (100 identities for CUHK03 and 751 for Market-1501) to 1,000 identities, *all* conventional supervised re-id models suffered severely, e.g. a 10-fold drop at Rank-1 for XQDA on CUHK03. More importantly, even though the supervised models were trained on a large-sized pre-labelled data in D_{p2} with an average of 3,483 cross-view images of 360 identities on CUHK03, and 7,737 images of 501 identities on Market-1501, their re-id performance was still significantly outperformed by *human-in-the-loop* models with 10-fold less human verification effort. This suggests the necessity of *human-in-the-loop* in real-world person re-id applications when the gallery population size becomes inevitably large. Moreover, to learn functionable supervised models, substantially more exhaustive pre-labelled training data are required. Such results suggest that *human-in-the-loop re-id* is a much better strategy for more efficiently exploiting human labour in real-world applications.

Second, HVIL improves significantly over the state-of-the-art human-in-the-loop model POP on Rank-1 score: from 46.3% to 56.1% on CUHK03 ($\sim 10\%$ in absolute terms) and from 56.0% to 78.0% on Market-1501 (over 20% in absolute terms). HVIL’s advantage continues over all ranks. This demonstrates compellingly the advantages of the HVIL model in cumulatively exploiting human verification feedback, whilst the existing human-in-the-loop models have no mechanisms for sharing human feedback knowledge among different probes.

Statistics analysis on human verification - Fig. 2 shows the comparisons of Rank-1 and Expected Rank (ER) on the 4 human-in-the-loop models over three verification feedback rounds. It is evident that the proposed HVIL model is more effective than the other three models in boosting Rank-1 scores and pushing up true matches’ ranking orders. The reasons are: (1) Given a large gallery population with potentially complex manifold structure, it is difficult to perform accurately graph label propagation for graph-based methods like POP and EMR. (2) Unlike POP/EMR/Rocchio, the proposed HVIL model optimises on re-id ranking losses (Eqn. (3)) specifically designed to maximise the three types of human verification feedback. (3) The HVIL model enables knowledge cumulation (Eqn. (5)). This is evident in Fig. 2 where HVIL yields notably better (lower) Expected Ranks (ER), even for the initial ER before verification feedback takes place on a probe (due to benefiting cumulative effect from other probes). In contrast, other models do not improve initial ER on each probe due to the lack of a mechanism to cumulate experience.

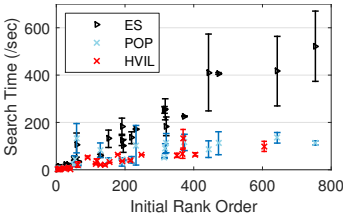


Fig. 3. Search time from different human-in-the-loop models on the same 25 randomly selected probes.

| Dataset | CUHK03 [9] | | | Market-1501 [35] | | |
|--------------------------------|-------------|------|------------|------------------|------|------------|
| Method | HVIL | POP | ES | HVIL | POP | ES |
| Found-matches(%) \uparrow | 56.1 | 46.3 | 100 | 78.0 | 56.0 | 100 |
| Browsed-images \downarrow | 32.9 | 42.1 | 234.1 | 19.5 | 38.3 | 108.7 |
| Feedback \downarrow | 2.1 | 2.3 | - | 1.6 | 1.9 | - |
| Search-time(sec.) \downarrow | 31.7 | 58.1 | 172.8 | 28.1 | 49.8 | 106.2 |

Table 2. Human verification effort vs. benefit. All measures are from averaging over all probes. \downarrow : lower better; \uparrow : higher better.

We further evaluated the human verification effort in relation to re-id performance benefit, collected from the human-in-the-loop re-id evaluation experiments reported above. We compared the HVIL model with the POP model and Exhaustive Search (ES) where a user performs exhaustive visual searching over the whole gallery ranking list (1,000) until finding a true match. The averaged statistics over all 6 trials were compared in Table 2. It is evident that though ES is guaranteed to locate a true match for every probe if it existed, it is much more expensive than POP (3 \times) and HVIL (5 \times) in search time given a 1,000-sized gallery. This difference will increase further on larger galleries. Comparing HVIL and POP, it is evident that HVIL is both more cost-effective (less Search-time, Browsed-images and Feedback) and more accurate (more Found-matches).

To better understand model convergence given human feedback, we conducted a separate experiment to measure the search time by different human-in-the-loop models given the initial rank lists on 25 randomly selected probes verified by multiple users. This experiment was evaluated by 10 independent sessions with the same set of 25 probes provided. In each session, the users were required to find a true match for all 25 probes. Specifically, for HVIL and POP, if a true match was not identified after 3 (maximum) feedback, the users then performed an exhaustive searching until it was found. The search time statistics for all 25 probes are shown in Fig. 3, where a bar shows the variance between 10 different sessions. It is unsurprising that ES is the least efficient whilst HVIL is the quickest in finding a true match, i.e. the data points of HVIL are much lower in search time. Moreover, it is evident that HVIL yields much better initial ranks, i.e. the data points of HVIL are more centred towards the bottom-left corner. This further shows the benefit of cumulative learning in HVIL (Sec. 2.3). Fig. 4 shows two visual examples of the HVIL model in action, where user feedback efficiently push true matches to top ranks within 2 rounds of interactions.

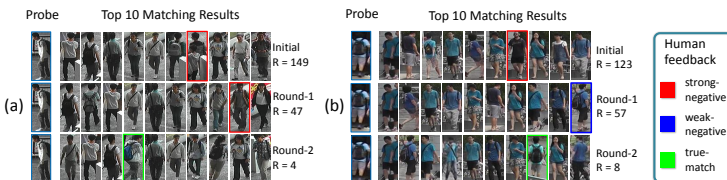


Fig. 4. HVIL re-id examples on CUHK03 (a) and Market-1501 (b). The ranks of true matches before user feedback are shown.

| Dataset | CUHK03 [9] ($N_g = 360$) | | | | Market-1501 [35] ($N_g = 501$) | | | | VIPeR ($N_g = 316$) | | | |
|------------------|----------------------------|-------------|-------------|-------------|----------------------------------|-------------|-------------|-------------|-----------------------|-------------|-------------|-------------|
| Rank (%) | 1 | 5 | 10 | 20 | 1 | 5 | 10 | 20 | 1 | 5 | 10 | 20 |
| L2 | 4.6 | 14.0 | 21.1 | 28.7 | 23.0 | 44.0 | 55.1 | 65.7 | 14.7 | 28.0 | 40.6 | 52.1 |
| kLFDA [10] | 6.2 | 19.0 | 28.3 | 39.1 | 29.1 | 58.9 | 71.2 | 82.2 | 32.3 | 65.8 | 79.7 | 90.9 |
| XQDA [11] | 5.3 | 14.2 | 21.1 | 30.0 | 28.7 | 54.5 | 65.6 | 75.3 | 40.0 | 68.1 | 80.5 | 91.1 |
| MLAPG [12] | 5.3 | 15.2 | 23.5 | 33.9 | 25.2 | 51.4 | 65.3 | 77.4 | 40.7 | 69.9 | 82.3 | 92.4 |
| HVIL - M_{avg} | 5.8 | 17.6 | 26.3 | 36.3 | 27.3 | 56.7 | 68.2 | 80.1 | 21.8 | 51.0 | 66.3 | 82.4 |
| HVIL - M_τ | 6.5 | 19.0 | 27.4 | 37.6 | 31.6 | 60.1 | 72.7 | 83.5 | 34.7 | 63.2 | 78.0 | 90.3 |
| HVIL - RMEL | 9.3 | 20.7 | 29.0 | 39.5 | 33.8 | 61.0 | 73.6 | 83.5 | 42.4 | 72.6 | 83.0 | 90.4 |

Table 3. Evaluating automated person re-id with CMC performances.

4.2 Evaluation on Automated Person Re-Id

The proposed RMEL model was evaluated for *automated person re-id* against both state-of-the-art supervised models and baseline ensemble models as follows.

Training/testing protocol - In each of the overall 6 trials, we employed the human verified true matches on D_{p1}^i (168 pairs on CUHK03 and 234 pairs on Market-1501 in average, as not all probe images found their true matches with a maximum of three feedback) to learn the weights for constructing a strong ensemble model using all the verified weak models $\{M_j\}_{j=1}^\tau$ collected from our previous experiments on *human-in-the-loop re-id*. The strong ensemble model was then deployed for testing on the separate partition D_{p2}^i with the gallery size of 360 and 501 for CUHK03 and Market-1501 respectively. For performance evaluation, we adopted the standard single-shot test setting, i.e. randomly sampling 360 cross-camera person image pairs from CUHK03 and 501 pairs from Market-1501 on $\{D_{p2}^i\}_{i=1}^6$ to construct the test gallery and probe sets over six trials. The averaged CMC performance over all trials was reported.

Competitors A - Three state-of-the-art supervised re-id models are compared: kLFDA [10], XQDA [11], and MLAPG [12] were trained using 300 ground-truth labelled data from \mathcal{P}^i (300) and \mathcal{G}^i (1,000) of D_{p1}^i , for both CUHK03 and Market-1501. The trained models were tested on the separate partition D_{p2}^i with same testing protocol as above.

Competitors B - For fully evaluating the effect of the HVIL-RMEL model, two more ensemble baseline models are compared: (1) HVIL - M_τ : The incrementally optimised re-id model M_τ obtained by HVIL from the last probe image at time τ during the *human-in-the-loop* process. (2) HVIL - M_{avg} : A naive approach to ensemble weak models, that is, simply taking an average weighting of all weak models $\{M_j\}_{j=1}^\tau$ as the ensemble re-id model.

Results on CUHK03 and Market-1501 - Table 3 reports the result. For CUHK03, there is insufficient labelled data for all camera pairs during training, given only one pair of randomly selected single-shot images per identity. All models including HVIL-RMEL generated poor re-id performances (Rank-1 < 10%), much less than state-of-the-art reported in the literature. For Market-1501, a similar problem exists although less pronounced. Note, the results in Table 3 are based on a single-shot test setting. This is a much harder problem than the multi-shot test setting [35] where on average 14.8 true matches exist in the gallery for each probe. When HVIL-RMEL was evaluated under the same multi-shot setting on Market-1501, it yields 53.5%, 83.0%, 89.0%, 94.1%

for Rank-1/5/10/20 respectively, significantly outperforms [35]. Given the experimental results above, it is evident that: Due to (1) a much larger unlabelled test gallery population than the labelled training set, (2) a lack of sufficient multi-shot training/testing data in many camera pairs, *human-in-the-loop* approach to re-id is not only desirable, but essential for re-id in real world applications.

Nevertheless, for *automated person re-id*, the proposed HVIL-RMEL still achieves the best performance among all models with a Rank-1 of 9.3% on CUHK03 and 33.8% on Market-1501. More importantly, even though less true-match data (168 pairs for CUHK03 and 234 pairs for Market-1501) were used to learn the ensemble weighting for the RMEL model as compared to the ground-truth data (300 pairs for both benchmarks) used to train kLFDA, XQDA and MLAPG, it is evident that the human verification feedback process yields more discriminative information for optimising probe re-id directly in the gallery population, resulting in a more optimal ensemble model. It is also evident that naively taking an average ensemble model (HVIL - \mathbf{M}_{avg}) gives even poorer performance than the cumulatively learned single model (HVIL - \mathbf{M}_τ).

Results on VIPeR - To compare HVIL-RMEL in a more comparable context defined in the literature on *automated person re-id*, we tested the HVIL-RMEL model on the VIPeR [13] benchmark under the exact setting of the established protocol: splitting the 632 identities into 50–50% partitions for training and testing sets. For obtaining weak re-id models, we simulated HVIL feedback update by simply giving the ground-true matching pairs instead of weak/strong-negatives (Eqn. (9)); therefore each weak model was obtained by a true-match, using the same information as training a conventional supervised model. The last/right panel of Table 3 compares the performance of such a HVIL-RMEL model against the published results of kLFDA, XQDA and MLAPG³. It is evident that the proposed model yields state-of-the-art performance under the same conventional re-id settings, with Rank-1 score of 42.4%, slightly better, by 2.4% and 1.7% respectively, than the current state-of-the-art XQDA and MLAPG.

5 Conclusions

We formulated a novel approach to human-in-the-loop person re-id by introducing a Human Verification Incremental Learning (HVIL) model, designed to overcome two unrealistic assumptions adopted by existing re-id models that prevent them to be scalable to real world applications. In particular, the proposed HVIL model avoids the need for collecting off-line pre-labelled training data and is scalable to re-id tasks in large gallery sizes. The advantage of HVIL over other human-in-the-loop models is its ability to learn cumulatively from human feedback on more probe images when available. We further developed a regularised metric ensemble learning (RMEL) method to explore HVIL for automated re-id tasks when human feedback is unavailable. Extensive comparisons on the CUHK03 [9] and the Market-1501 [35] benchmarks show the potentials of the proposed HVIL-RMEL model for real-world re-id deployments.

³ A different 26,960-dim LOMO feature [11] were used for the published XQDA and MLAPG results [12, 11] shown in Table 3. They were worsened using the 5,138-dim feature [47] adopted in our experiments, not shown here due to space limitation.

References

1. Gong, S., Cristani, M., Yan, S., Loy, C.C.: Person re-identification. Volume 1. Springer (2014)
2. Mignon, A., Jurie, F.: Pcca: A new approach for distance learning from sparse pairwise constraints. In: IEEE Conference on Computer Vision and Pattern Recognition, Providence, Rhode Island, United States (June 2012)
3. Koestinger, M., Hirzer, M., Wohlhart, P., Roth, P.M., Bischof, H.: Large scale metric learning from equivalence constraints. In: IEEE Conference on Computer Vision and Pattern Recognition, Providence, Rhode Island, United States (June 2012)
4. Zheng, W.S., Gong, S., Xiang, T.: Re-identification by relative distance comparison. IEEE Transactions on Pattern Analysis and Machine Intelligence (March 2013) 653–668
5. Pedagadi, S., Orwell, J., Velastin, S.A., Boghossian, B.A.: Local fisher discriminant analysis for pedestrian re-identification. In: IEEE Conference on Computer Vision and Pattern Recognition, Portland, Oregon, United States (June 2013)
6. Zhao, R., Ouyang, W., Wang, X.: Learning mid-level filters for person re-identification. In: IEEE Conference on Computer Vision and Pattern Recognition, Columbus, Ohio, United States (June 2014)
7. Wang, T., Gong, S., Zhu, X., Wang, S.: Person re-identification by video ranking. In: European Conference on Computer Vision, Zurich, Switzerland (September 2014)
8. Wang, T., Gong, S., Zhu, X., Wang, S.: Person re-identification by discriminative selection in video ranking. IEEE Transactions on Pattern Analysis and Machine Intelligence (January 2016)
9. Li, W., Zhao, R., Xiao, T., Wang, X.: Deepreid: Deep filter pairing neural network for person re-identification. In: IEEE Conference on Computer Vision and Pattern Recognition, Columbus, Ohio, United States (June 2014)
10. Xiong, F., Gou, M., Camps, O., Sznaiar, M.: Person re-identification using kernel-based metric learning methods. In: European Conference on Computer Vision, Zurich, Switzerland (September 2014)
11. Liao, S., Hu, Y., Zhu, X., Li, S.Z.: Person re-identification by local maximal occurrence representation and metric learning. In: IEEE Conference on Computer Vision and Pattern Recognition, Boston, Massachusetts, United States (June 2015) 2197–2206
12. Liao, S., Li, S.Z.: Efficient psd constrained asymmetric metric learning for person re-identification. In: IEEE International Conference on Computer Vision. (December 2015)
13. Gray, D., Brennan, S., Tao, H.: Evaluating appearance models for recognition, reacquisition and tracking. In: IEEE International Workshop on Performance Evaluation for Tracking and Surveillance. (2007)
14. Zheng, W.S., Gong, S., Xiang, T.: Associating groups of people. In: British Machine Vision Conference. (2009)
15. Liu, C., Loy, C.C., Gong, S., Wang, G.: Pop: Person re-identification post-rank optimisation. In: IEEE International Conference on Computer Vision, Sydney, Australia (December 2013)
16. Schapire, R.E.: The strength of weak learnability. Machine Learning 5(2) (July 1990) 197–227

17. Amit, Y., Geman, D.: Shape quantization and recognition with randomized trees. *Neural Comput.* **9**(7) (October 1997) 1545–1588
18. Gong, S., Cristani, M., Chen, C.L., Hospedales, T.M.: The re-identification challenge. In: *Person Re-Identification*. Springer (2014)
19. Ding, S., Lin, L., Wang, G., Chao, H.: Deep feature learning with relative distance comparison for person re-identification. *Pattern Recognition* (2015) 2993–3003
20. Paisitkriangkrai, S., Shen, C., van den Hengel, A.: Learning to rank in person re-identification with metric ensembles. In: *IEEE Conference on Computer Vision and Pattern Recognition*. (2015) 1846–1855
21. Ustinova, E., Ganin, Y., Lempitsky, V.: Multiregion bilinear convolutional neural networks for person re-identification. arXiv preprint arXiv:1512.05300 (2015)
22. Shi, H., Zhu, X., Liao, S., Lei, Z., Yang, Y., Li, S.Z.: Constrained deep metric learning for person re-identification. arXiv preprint arXiv:1511.07545 (2015)
23. Ahmed, E., Jones, M.J., Marks, T.K.: An improved deep learning architecture for person re-identification. In: *CVPR, IEEE* (2015) 3908–3916
24. Liu, X., Song, M., Tao, D., Zhou, X., Chen, C., Bu, J.: Semi-supervised coupled dictionary learning for person re-identification. In: *IEEE Conference on Computer Vision and Pattern Recognition*, Columbus, Ohio, United States (June 2014)
25. Kodirov, E., Xiang, T., Gong, S.: Dictionary learning with iterative laplacian regularisation for unsupervised person re-identification. In: *British Machine Vision Conference*, Swansea, United Kingdom (September 2015)
26. Farenzena, M., Bazzani, L., Perina, A., Murino, V., Cristani, M.: Person re-identification by symmetry-driven accumulation of local features. In: *IEEE Conference on Computer Vision and Pattern Recognition*, San Francisco, California, United States (June 2010)
27. Zhao, R., Ouyang, W., Wang, X.: Unsupervised salience learning for person re-identification. In: *IEEE Conference on Computer Vision and Pattern Recognition*, Portland, Oregon, United States (June 2013)
28. Wang, H., Gong, S., Xiang, T.: Unsupervised learning of generative topic saliency for person re-identification. In: *British Machine Vision Conference*, Nottingham, United Kingdom (September 2014)
29. Layne, R., Hospedales, T.M., Gong, S.: Domain transfer for person re-identification. In: *Workshop of ACM International Conference on Multimedia*, Barcelona, Catalunya, Spain (October 2013)
30. Ma, A.J., Yuen, P.C., Li, J.: Domain transfer support vector ranking for person re-identification without target camera label information. In: *IEEE International Conference on Computer Vision*, Sydney, Australia (December 2013)
31. Wang, X., Zheng, W.S., Li, X., Zhang, J.: Cross-scenario transfer person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology* **PP**(99) (June 2015)
32. Ma, A.J., Li, J., Yuen, P.C., Li, P.: Cross-domain person reidentification using domain adaptation ranking svms. *IEEE Transactions on Image Processing* **24**(5) (2015) 1599–1613
33. Das, A., Panda, R., Roy-Chowdhury, A.: Active image pair selection for continuous person re-identification. In: *IEEE International Conference on Image Processing*, Quebec, Canada (September 2015)
34. Hirzer, M., Beleznai, C., Roth, P.M., Bischof, H.: Person re-identification by descriptive and discriminative classification. In: *Proceedings of the 17th Scandinavian Conference on Image Analysis*. SCIA'11, Berlin, Heidelberg, Springer-Verlag (2011) 91–102

35. Zheng, L., Shen, L., Tian, L., Wang, S., Wang, J., Tian, Q.: Scalable person re-identification: A benchmark. In: Proceedings of the IEEE International Conference on Computer Vision. (2015) 1116–1124
36. Lim, D., Lanckriet, G.: Efficient learning of mahalanobis metrics for ranking. In Jebara, T., Xing, E.P., eds.: International Conference on Machine learning. (2014) 1980–1988
37. Weston, J., Bengio, S., Usunier, N.: Large scale image annotation: Learning to rank with joint word-image embeddings. In: European Conference of Machine Learning. (2010)
38. Usunier, N., Buffoni, D., Gallinari, P.: Ranking with ordered weighted pairwise classification. In: Proceedings of the 26th Annual International Conference on Machine Learning. ICML '09, New York, NY, USA, ACM (2009) 1057–1064
39. Chechik, G., Sharma, V., Shalit, U., Bengio, S.: Large scale online learning of image similarity through ranking. *J. Mach. Learn. Res.* (March 2010) 1109–1135
40. Bottou, L.: Large-scale machine learning with stochastic gradient descent. In: Proceedings of COMPSTAT'2010. Springer (2010) 177–186
41. Tsuda, K., Rätsch, G., Warmuth, M.K.: Matrix exponentiated gradient updates for on-line learning and bregman projection. In: *Journal of Machine Learning Research.* (2005) 995–1018
42. Higham, N.J.: Matrix nearness problems and applications. University of Manchester. Department of Mathematics (1988)
43. Kivinen, J., Warmuth, M.K.: Exponentiated gradient versus gradient descent for linear predictors. *Information and Computation* (1997) 1–63
44. Jain, P., Kulis, B., Dhillon, I.S., Grauman, K.: Online metric learning and fast similarity search. In: *Advances in Neural Information Processing Systems*, Vancouver, British Columbia, Canada (December 2009) 761–768
45. Grant, M., Boyd, S.: CVX: Matlab software for disciplined convex programming, version 2.1. <http://cvxr.com/cvx> (March 2014)
46. Boyd, S., Vandenberghe, L.: *Convex Optimization*. Cambridge University Press, New York, NY, USA (2004)
47. Lisanti, G., Masi, I., Del Bimbo, A.: Matching people across camera views using kernel canonical correlation analysis. In: *ACM International Conference on Distributed Smart Cameras*, Venice, Italy (November 2014)
48. Wang, X., Han, T.X., Yan, S.: An hog-lbp human detector with partial occlusion handling. In: *IEEE International Conference on Computer Vision*, Kyoto, Japan (September 2009)
49. Ahonen, T., Hadid, A., Pietikainen, M.: Face description with local binary patterns: Application to face recognition. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* (2006) 2037–2041
50. Lin, W.C., Chen, Z.Y., Ke, S.W., Tsai, C.F., Lin, W.Y.: The effect of low-level image features on pseudo relevance feedback. *Neurocomputing* (October 2015) 26–37
51. Datta, R., Joshi, D., Li, J., Wang, J.Z.: Image retrieval: Ideas, influences, and trends of the new age. *ACM Computing Surveys* (April 2008) 5:1–5:60
52. Xu, B., Bu, J., Chen, C., Cai, D., He, X., Liu, W., Luo, J.: Efficient manifold ranking for image retrieval. In: *ACM SIGIR Conference on Research and Development in Information Retrieval*, Beijing, China (July 2011) 525–534