



Contents lists available at ScienceDirect

Transportation Research Part C

journal homepage: www.elsevier.com/locate/trc



Deep visual Re-identification with confidence

George Adaimi ^{*}, Sven Kreiss, Alexandre Alahi

VITA, EPFL, Switzerland

ARTICLE INFO

Keywords:

Traffic monitoring
Person re-identification
Vehicle re-identification
Flow monitoring

ABSTRACT

Transportation systems often rely on understanding the flow of vehicles or pedestrian. From traffic monitoring at the city scale, to commuters in train terminals, recent progress in sensing technology make it possible to use cameras to better understand the demand, *i.e.*, better track moving agents (*e.g.*, vehicles and pedestrians). Whether the cameras are mounted on drones, vehicles, or fixed in the built environments, they inevitably remain scatter. We need to develop the technology to re-identify the same agents across images captured from non-overlapping field-of-views, referred to as the visual re-identification task. State-of-the-art methods learn a neural network based representation trained with the cross-entropy loss function. We argue that such loss function is not suited for the visual re-identification task hence propose to model confidence in the representation learning framework. We show the impact of our confidence-based learning framework with three methods: label smoothing, confidence penalty, and deep variational information bottleneck. They all show a boost in performance validating our claim. Our contribution is generic to any agent of interest, *i.e.*, vehicles or pedestrians, and outperform highly specialized state-of-the-art methods across 6 datasets. The source code and models are shared towards an open science mission.

1. Introduction

An important goal of transportation research is to improve and provide efficient public transportation systems that can accommodate many agents, whether it be vehicles or pedestrians, every day. This is especially important nowadays with the huge traffic congestion costing billions of dollars (Schneider, 2018). As a result, research efforts have been directed towards management and control of vehicle and pedestrian flows. Important prerequisites for such transportation network analysis are origin–destination (OD) matrices, which allow researchers to understand a population's trip demand. With the recent developments in new methods, such as data-driven methods (Krishnakumari et al., 2019), OD estimation have achieved good performance in various tasks of traffic management. However, it still faces the issue of how representative the chosen samples are of the population. One way to deal with this problem is to collect more data from the population, which is expensive and time-consuming when using traditional methods such as surveys or interviews. Another way is to make use of smartphones to collect such data (Zhao et al., 2015). Visual re-identification is a faster and cheaper way to collect these data. Visual re-identification is the task of associating images of the same agent taken from different cameras or from the same camera in different occasions. This is represented in Fig. 1. This task is nowadays possible due to the complex network of cameras already places in and around cities. Moreover, recent works have been pushing towards the use of drone technology to collect massive amounts of data in order to study the traffic phenomena, such as the recently released pNEUMA dataset (Barmpounakis and Geroliminis, 2020). All these collected data can be used as inputs to a visual re-identification model to associate different agents together and obtain their paths.

* Corresponding author.

E-mail addresses: george.adaimi@epfl.ch (G. Adaimi), sven.kreiss@epfl.ch (S. Kreiss), alexandre.alahi@epfl.ch (A. Alahi).

<https://doi.org/10.1016/j.trc.2021.103067>

Received 5 June 2020; Received in revised form 15 January 2021; Accepted 23 February 2021

Available online 22 March 2021

0968-090X/© 2021 The Authors. Published by Elsevier Ltd. This is an open access article under the CC BY license

(<http://creativecommons.org/licenses/by/4.0/>).



Fig. 1. In this work, we present a new visual re-identification method, *i.e.*, whether agents (vehicles or pedestrians) captured in different images in non-overlapping views belong to the same agent. Agents with the orange bounding box belong to the same agent. Note different agents can be visually very similar.

The task of re-identification (re-id) has long been a task of extracting features/representations from two observations and measuring how similar these features are. Since different variations affect these features, many works have introduced different methods to improve their extraction. Initially, these features were hand-crafted and include spatio-temporal information such as color, width and height, and salient edge histograms. Some work have also tried to use different input modalities such as depth (Wu et al., 2017a; Karianakis et al., 2017; Wu et al., 2015), infrared (Wu et al., 2017b), LiDAR (Zhao et al., 2019), or Inductive Loop Detectors for vehicles (Abdulhai and Tabib, 2003). These features, however, fail drastically when dealing with unexpected scenarios. To remedy this problem and with the advent of machine learning, researchers are now benefiting from the strength of deep learning to be able to extract more general and more discriminative features allowing them to reach high performance. Since then, an arms race of methods was built on top of this by making use of different object-specific characteristics (*e.g.*, human semantic segmentation, pose) and by learning features through the supervision of a cross-entropy loss.

A main pitfall of learning with cross-entropy supervision is the fact that it separates the different inputs solely based on the labels without taking into consideration the actual similarity between the inputs. None of the recent methods have tackled the problem where, even though two very similar agents are distinct, their similarity score should encode information about how similar they appear while also distinguishing them. The network usually tries its best to find a boundary between the different classes even for inputs that are very similar. This leads the network to find unreasonable explanations for the differences in labels and thus would negatively affect its generalizability. Since re-id also deals with the problem of having a small set of images per class, it would aggravate this issue. Controlling the network's confidence in its predictions would alleviate this problem. To the best of your knowledge, this is the first paper to apply this concept in the field of re-identification.

In this paper, we propose to model confidence when learning representations appropriate for re-id. By inducing doubt while training a network, we are able to tackle the inherent problem discussed previously when cross-entropy is used in a distance metric and representation learning problem such as person or vehicle re-identification. Inspired by previous works that use uncertainty to regularize the network, we study three alternatives that aim at reducing the confidence of the network and show a gain in mAP and Rank-1 across 6 different datasets. Although these methods have shown only a small improvement in other image classification tasks (Szegedy et al., 2015; Alemi et al., 2016; Pereyra et al., 2017), they drastically improve the performance of re-id models due to its innate problem (Section 3). By combining our methods with advanced ranking methods, we outperform state-of-the-art models

without modeling additional characteristics specific to the object in question. The software is open-source and is made available online.¹

2. Related work

Initially, re-identification's origins were based on multi-camera tracking (Alahi et al., 2017a). In this context, both visual and spatio-temporal information were used to predict whether the same agent was found in different cameras and at different times (Wang, 2013; Huang and Russell, 1997; Alahi et al., 2017b; Mazzon et al., 2012; Held et al., 2014). Gheissari et al. (2006) were the first to make use of only visual information to match different pedestrians. They made use of color and salient edge histograms to re-identify around 44 pedestrians recorded by 3 different cameras. This work clearly divided visual re-identification from multi-camera tracking. Alahi et al. (2014, 2008) showed the benefits of visual re-identification in a network of fixed and mobile cameras. Fixed cameras installed at urban intersection can help the detection algorithms of cameras mounted on vehicles. Since then, researchers have proposed many methods to solve the visual re-identification tasks (Corvee et al., 2010b,a). In the remaining of this section, we briefly present key methods tackling the problem for “person” and “vehicle”.

2.1. Person re-identification

Even though person re-identification is highly beneficial to analyze how people make use of different transportation modes, it is not a widely studied matter in this context. Recent works in the transportation field have been using traditional and linear methods such as the Kalman filter (Zhao et al., 2019; Guo et al., 2016) and particle filters (Gaddigoudar et al., 2017; Li et al., 2016; Hue et al., 2002) to track pedestrians by using spatio-temporal information. These methods work well when tracking for a short period and while assuming that pedestrian's movement is linear. Applying the same methods over multiple cameras as well as in more complicated environments such as in a train station would lead to poor identification of passengers.

With the prevalence of deep neural networks in most computer vision tasks, person re-id followed this success when Li et al. (2014) introduced a deep learning method for re-id that tried to overcome the problems of bounding box misalignment, photometric and geometric transforms while also introducing a new bigger dataset specifically for this task. This paved the way for new methods and datasets to emerge, causing the person re-id performance of machines to improve. Other work developed new methods that tackle specific challenges in person re-id by introducing different architectures and modules (Li et al., 2017a, 2018; Song et al., 2018; Wang et al., 2018a).

Attention in Person reID. Recently, many methods have tried to improve the representation of the input by training multiple networks that extract global and local features and then combining these features to form the final representation. This is usually done by using either a deterministic way of dividing the different parts of the representation (Yi et al., 2014; Cheng et al., 2016; Varior et al., 2016; Ahmed et al., 2015) or making use of attention modules to separate the different parts (Li et al., 2018; Zhao et al., 2017; Yao et al., 2017). Other works extracted intermediate representations to gain information about the input at different levels arguing that this allows the network to learn distinctive characteristics of the input at different scales (Chang et al., 2018; Wang et al., 2018b,a). Even though these methods showed improvement over their predecessors, they usually require separate networks to process each of the different features, leading to a more complex architecture and training procedure.

Human Characteristics. Another direction other researchers have taken is to make use of information and characteristics related specifically to humans in order to improve person re-id. The work by Xu et al. (2018) aims at detecting three different types of pose information such as keypoints, rigid body parts (e.g., torso), and non-rigid parts (e.g., limbs). These information were extracted using an off-the-shelf human pose estimator. Then, with the help of these body parts, the features extracted by a feature extractor are refined and used to classify the different pedestrians. The use of third-party methods makes their model highly dependent on the performance of these methods. Another approach by Saquib Sarfraz et al. (2018) uses keypoint information, in addition to the input image, to train a ResNet-50 model as well as another connected module that detects the view (front, back or side). Kalayeh et al. (2018) also made use of features extracted from different body parts and concatenated them to form a global feature which in turn was used to perform re-identification. The disadvantage of these methods is their high dependence on other methods and datasets that require annotation. Moreover, the fact that these models depend on specific human characteristics prevents them from being leveraged for other image-retrieval and clustering tasks.

Re-Ranking. In addition to learning better features, many works have tried to improve the ranking process of person re-id by including information about how the different galleries are related instead of just using the relationship between the pairs of queries and galleries (Zhong et al., 2017a; Ye et al., 2015; Shen et al., 2018b; Zhong et al., 2017a; García et al., 2015; Leng et al., 2015; Ye et al., 2015; Ye et al., 2016). Zhong et al. (2017a) introduced a method for refining the distances between the queries and galleries by making use of the k-reciprocal nearest neighbors. This is done as a post-processing step to improve the ranking process. Shen et al. (2018b) argued that this does not help in learning better features during training and introduced a new learnable module that performs a random walk on a graph connecting the different gallery images. By performing a random walk operation, gallery-to-gallery (G2G) information is taken into consideration while training the network, thus resulting in a

¹ <https://git.io/deep-visual-reID-confidence>.

more complete representation that provides a better ranking performance. Other methods also tried to include G2G information by using Graph Neural Networks (Shen et al., 2018a) and Conditional Random Fields (Chen et al., 2018). We will make use of G2G information by applying different re-ranking methods.

Metric Learning. Several previous works have tried to tackle the problem of person re-id by introducing new metric loss functions. Both contrastive (Hadsell et al., 2006) and binary loss functions have been employed in order to push apart negative image pairs while pulling positive image pairs together (Rama Varior et al., 2016; Ahmed et al., 2015). Taking into consideration both the pull and push of contrastive loss, other methods (Wang et al., 2018a; Liu et al., 2018; Wang et al., 2018b) used triplet loss that simultaneously tackles negative and positive pairs leading to a less greedy method. Chen et al. (2017) extended this loss to quadruplet inputs. The drawback of these methods is their high sensitivity to the sampling technique used. As a result, Yu et al. (2018) introduced the HAP2S loss to tackle this drawback and showed improvement in performance. All the above methods try to encode metric information in the embedding space compared to cross-entropy which is considered as a representation learning method.

2.2. Vehicle re-identification

Vision-based methods for vehicle re-identification are rather new. Initially, vehicle datasets were small and mainly used for car color, model classification, or detection (Tian et al., 2015). As a result, Liu et al. (2016a) built a large scale dataset for vehicle re-identification similar to person re-id datasets. They also made use of many techniques derived for person re-id and compared their performance. They later extended their work and added license plate verification as a way to improve re-identification (Liu et al., 2016b). This intuition comes from the fact that each car has a unique license plate. Metric learning techniques were also extended for this task. Liu et al. (2016) introduced coupled clusters loss, a variant of triplet loss, to deal with sensitivity to the choice of the triplet samples. Zhang et al. (2017) also improved triplet loss by augmenting the training with a classification loss and modifying the dataset sampling method. Instead of randomly sampling triplets of anchors, positive, and negative samples, their method ensure that the negative sample is an anchor or positive sample in another triplet. This provides a way for negatives to be pushed towards similar images rather than being pushed away from the anchor randomly. Bai et al. (2018) tried to deal with the problem of inter-class similarity and intra-class variance by introducing the group sensitive triplet embedding (GSTE). This is done by combining samples into intermediate “groups” at different granularity levels such as vehicle ID and vehicle model.

Other works made use of different attributes and modalities to improve vehicle re-identification. Li et al. (2017b) performed multi-task training that includes ID classification, attribute recognition, contrastive loss, and triplet loss. Tang et al. (2017) made use of hand-crafted features such as color and introduced a multi-modal metric learning method to fuse these features with deep features extracted using a neural network. Moreover, GANs are being increasingly adopted in vehicle re-identification. The main intuition is to transfer the query image to a domain that makes it more robust and efficient to compare with other images (Zhou and Shao, 2017; Lou et al., 2019a). Zhou and Shao (2017) proposed the generation of vehicle images from different views to deal with cross-view re-identification.

Large-scale datasets for vehicle re-identification are still recent and with the emergence of intelligent transportation, more research is being developed to improve this field. While it does introduce certain challenges different from person re-identification, the ability to find a re-identification method that performs well on any object of interest is important. In this paper, we do not make use of characteristics specific to the object of interest or feature division and show the importance of confidence when training a re-id model with a cross-entropy loss.

3. Problem formulation

A re-id model’s main task is to distinguish between different agents across images. As previously stated, this is a challenging task since it tries to relate images of agents across different cameras as well as at different times. The fact that the images are captured under different circumstances might lead to subtle differences in hue and image color that can drastically effect the performance of a re-id model. Moreover, the illumination, background clutter, occlusion, and observable object parts are usually dramatically different which might easily fool the network and render it unusable. Even images captured by the same camera can have many of these variations.

Due to the challenges explained above, there is not always a clear margin of separation between individual agents. Pedestrians or vehicles in some cases have very subtle differences that separate them from each other making the task of identifying them even more challenging for an observer. A good example is shown in Fig. 2 which introduces the inherent challenge we are trying to tackle in this paper.

The pedestrians within the images in Fig. 2 are very difficult to discern from one another even for a human eye. Each pair of images show two different pedestrians who share very similar appearances. When a model is trained to separate these images, it might face difficulties doing so. Since the images are very similar and a network’s only main goal is to reduce its loss, it will learn to focus on the pose or even the illumination of the images to discern them. These two variations are some of the many variations that previous methods try to overcome. This problem is also shared with vehicle re-identification. Since different vehicles might share the same model and color, many variations such as illumination, viewing angle, and weather might be used by the network to separate them.

Current state-of-the-art re-id systems train their own models by using the cross-entropy loss function. The cross-entropy calculates the number of bits needed for an event, which in this case is the label given the input, using the estimated probability distribution



Fig. 2. Pairs of images of different IDs but very similar appearance. None of these images belong to the same agent. Images taken from Market1501 ([Zheng et al., 2015](#)) & VERI-Wild ([Lou et al., 2019b](#)) datasets.

instead of the true distribution. In the case of training a neural network, the cross-entropy is minimized so that the model distribution is the same as that of the ground-truth, which is usually a one-hot encoding. This means minimizing this loss pushes the distribution of the model to output a high probability for the correct label while outputting very low probabilities for the others. The fact that cross-entropy requires that the logits for the ground-truth label to be much bigger than other labels pushes the network to take into consideration certain destructive variations to separate the different classes and especially for images such as in Fig. 2.

In order to modify the cross entropy in a way that solves the problem described above, we add a missing term to the loss function which allows it to not be confident about certain data points. Thus, the modified cross entropy loss function allows the network not to overfit on variations that are destructive for the re-id task and accept the fact that pedestrians or vehicles do sometimes look very similar. The idea of preventing the network from being very confident is not a new concept. However, its evaluation on other computer vision tasks, such as object detection, only leads to slight improvements in performance. From the reasoning based on Fig. 2 as well as the characteristics of person and vehicle re-id datasets, we show in this paper that this concept, if applied to a simple baseline, can improve the results drastically and even outperform certain highly specialized state-of-the-art methods.

4. Method

Current re-id models face difficulties in distinguishing between different agents who share some visual similarities due to the model's objective of maximizing its confidence in its predictions, as previously discussed. In this section, we place the three common methods into a common framework of a cross-entropy term and a KL divergence term. This will allow us to investigate common properties and to identify their differences. We show how the three methods, LS, CP, and VIB, are specific instances of our common framework which forces the model to be less confident of its predictions than a plain cross entropy term. These methods usually show a small improvement when used during training in other computer vision tasks ([Szegedy et al., 2016](#); [Pereyra et al., 2017](#); [Alemi et al., 2016](#)); however, we show that, because of the problems specified in Section 3, these methods provide a drastic boost for the task of person and vehicle re-id.

We will bring all methods into the following common framework of a cross-entropy regularized with a KL divergence where the loss L is

$$L = \alpha H(q, p) + \beta KL \quad (1)$$

where $q = q(y|x)$ is the ground truth distribution of the output id y given the input image x , $p = p(y|x)$ is the predicted distribution, $H(q, p)$ is the cross-entropy between q and p and KL is the Kulback–Leibler divergence (Kullback and Leibler, 1951). In ReID, we have multiple ids $y \in Y$ where we denote the number of Y – the number of classes – as C .

4.1. Label smoothing

With Label Smoothing (Szegedy et al., 2016), the predicted distribution p is regularized towards the uniform distribution u with the KL divergence term:

$$L_{LS} = \alpha H(q, p) + \beta KL[u, p] \quad . \quad (2)$$

In this form, the KL divergence is the expectation over the uniform distribution of the logarithmic difference between u and p . Forming the convex mixture with $\alpha \equiv 1 - \beta$ and expanding this equation yields:

$$L_{LS} = -(1 - \beta) \sum_{y \in Y} q \log p + \beta \sum_{y \in Y} u \log \frac{u}{p} \quad (3)$$

$$= - \sum_{y \in Y} [\beta u + (1 - \beta)q] \log p \quad (4)$$

$$= H(q_{LS}, p) \quad (5)$$

where we dropped the constant term $u \log u$. We arrived at a single cross entropy without a KL divergence and defined a new label-smoothed ground truth distribution q_{LS} . The uniform distribution u is over the C classes, i.e., it evaluates to $1/C$. For a sample n :

$$q_{LS}(y|x_n) = \begin{cases} 1 - \frac{(C-1)\beta}{C} & \text{true } y, \\ \frac{\beta}{C} & \text{otherwise.} \end{cases} \quad (6)$$

From a KL divergence between u and p , we obtained a loss function similar to the regularizer introduced by Szegedy et al. (2016), which aims at allowing a model to be less confident about its prediction. It regularizes a softmax classifier by assigning a small value to all ground-truth labels.

This method makes sure that the label for the correct class does not become much larger than all other classes and thus prevents the network from over-fitting. When label smoothing was proposed and tested on ImageNet, it showed a small improvement of around 0.2% for top-1 error. Even though it did not show a huge improvement, we show in Section 7 that this method has a bigger effect on the task at hand based on the arguments stated in Section 3.

4.2. Confidence penalty

Reversing the arguments of the KL divergence, we obtain an equation for confidence penalty (Pereyra et al., 2017):

$$L_{CP} = \alpha H(q, p) + \beta KL[p, u] \quad (7)$$

Comparing Eqs. (7) to (2), we can observe the main difference. The error calculation, in this case, between the uniform and predicted distribution is weighted by the predicted distribution. Expanding this equation:

$$L_{CP} = \alpha H(q, p) + \beta \sum_{y \in Y} p \log \frac{p}{u} \quad (8)$$

$$= \alpha H(q, p) - \beta H(p) \quad (9)$$

where $\sum_{y \in Y} p \log u$ is removed since it is a constant. We notice that the resulting loss function aims at maximizing the entropy of the predicted distribution. The increase in entropy forces the network to be less certain of its predictions. Pereyra et al. (2017) reached a similar conclusion and showed that, by applying this method, they got a smoother predicted distribution as well as a small improvement on MNIST. This method, however, did not show an improvement on a more difficult dataset such as CIFAR-10. Similar to label smoothing, this method does not require any architecture modification as shown in Fig. 3.

4.3. Deep variational information bottleneck

Another way to increase the entropy of the output distribution is to force latent representations to be similar to each other irrespective of the input. This means that information distinguishing different samples is being lost. This idea can be derived from the information bottleneck (IB) principle (Tishby et al., 1999) where the mutual information between the input and latent representations is minimized. The IB principle (Tishby et al., 1999) is a technique that tries to find the best trade-off between

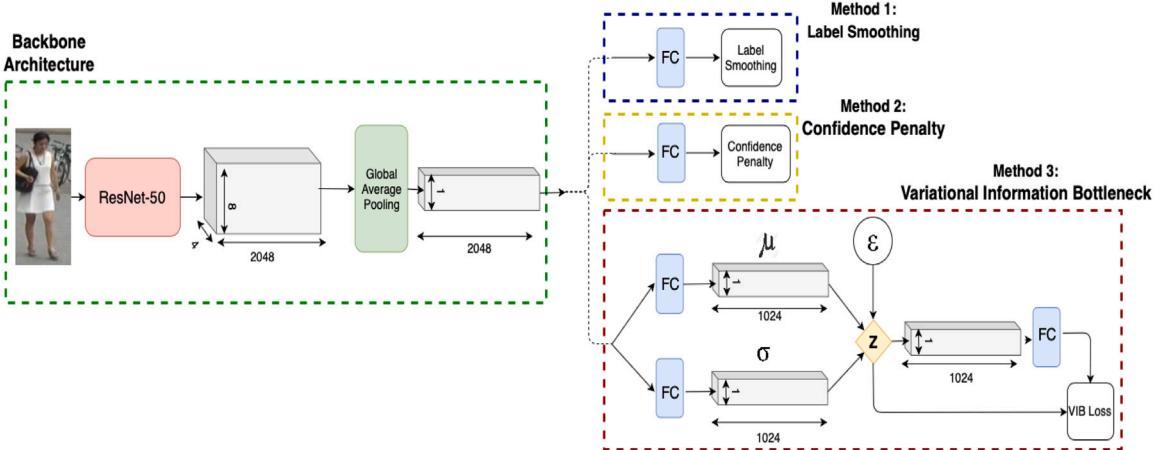


Fig. 3. Network architecture including the three methods being studied: Label Smoothing (LS), Confidence Penalty (CP), Variational Information Bottleneck (VIB). ϵ is the Gaussian noise needed for the reparameterization trick, μ and σ are the mean and standard deviation respectively of the latent Z distribution.

accuracy and complexity of latent variables. Latent variables are hidden variables that describe a specific input while maintaining all the relevant information needed for a specific task. The information bottleneck method tries to maximize this objective:

$$\max_{p(z|x)} \alpha I(z; y) - \beta I(z; x) \quad (10)$$

where z is the latent variable. Based on the above equation, the objective is to learn a representation z that is very informative about y while compressive about x . In order to apply the IB objective to a neural network, Alemi et al. (2016) approximated a lower bound to the information bottleneck by using variational inference and the reparameterization trick introduced by Kingma and Welling (2014) to introduce a new objective function referred to as Variational Information Bottleneck (VIB).

When applying this method, the model is divided into an encoder that takes the input x and maps it to a distribution describing the latent space z . The encoder outputs both the mean μ and standard deviation σ that describe this distribution. Then the predicted latent distribution is used to sample a specific latent representation. To force the second part of Eq. (10) to be maximized, this distribution should not depend on the input thus forcing the representation z to forget some information about it. This is done by minimizing the divergence between the encoder's distribution $w = w(z|x)$ and the prior r which is an isotropic multivariate Gaussian $r(z)$. The resulting objective function to minimize is:

$$L_{VIB} = \alpha H(q, \tilde{p}) + \beta KL[w, r] \quad (11)$$

where $\tilde{p} = p(y|f(x, \epsilon))$.

In order to compute the KL divergence analytically and back-propagate using its gradients, w is approximated by a multivariate Gaussian distribution with a diagonal covariance matrix. As can be seen in Eq. (11), if $\beta \rightarrow \infty$, the latent representation would follow a distribution independent of the input. This is somewhat similar to the effect of both confidence penalty and label smoothing where a single representation is forced to contain some information about more than one label. However, VIB applies this restriction directly to the latent space. Using this method while training, Alemi et al. (2016) showed close results to state-of-the-art models while using less information about the input which is measured using mutual information $I(x; z)$.

Compared to previously mentioned methods, in order to use the VIB loss, a fully connected layer is added at the output of the ResNet-50 base model to compute the mean and standard deviation as shown in Fig. 3.

5. Discussion

As can be seen in the equations above, all three methods try to increase the uncertainty of the model or in other words, decrease its confidence. On one hand, label smoothing and confidence penalty act on the output distribution while on the other hand, VIB acts on the latent representation directly. Thus this requires the original architecture to be modified to accommodate the VIB loss.

In addition, when expressed in terms of KL divergence, both label smoothing and confidence penalty are very similar except for the fact that the KL divergences are reversed. The forward KL divergence uses a constant represented by $u = \frac{1}{C}$ (Eq. (2)) to weight the log expression. However, the reverse KL divergence weighs using the output prediction of the model $p(y|x)$ which varies during the training process (Eq. (7)). In other words, confidence penalty does not equally penalize all the label predictions but implicitly gives more importance to predictions that the network incorrectly gives higher probabilities to and which are farther away than the uniform distribution. This weighing is adaptively changing as the network trains. Label smoothing however penalizes all label predictions equally. Moreover, label smoothing tries to prevent an output prediction of 0. This is because it is weighted by the uniform distribution u which is always greater than zero. Confidence penalty, however, might force certain output predictions to be zero although u is never zero.



Fig. 4. An example from the EPFL Roundabout dataset and the vehicles that are detected. We make use of the crop of vehicles to build a re-identification dataset similar to VERI-Wild.

By expressing all three methods in terms of KL divergence, we get more insight on why confidence penalty is able to outperform other methods. Compared to label smoothing, confidence penalty weighs the error by the network's current prediction and thus is more adaptable to the different inputs. Moreover, VIB provides the network with less degree of freedom compared to confidence penalty. This is because the latter acts on the output distribution allowing the representation, which is the main feature used for ranking in re-id, to move more freely in the feature space.

6. Experiments

To evaluate our proposed method, we use publicly available person and vehicle re-identification datasets which are Market-1501 ([Zheng et al., 2015](#)), MSMT17 ([Wei et al., 2018](#)), DukeMTMC-reID ([Ristani et al., 2016](#)), and VERI-Wild ([Lou et al., 2019b](#)).

Market1501 ([Zheng et al., 2015](#)): The Market dataset is a well-known person re-identification dataset that contains 32,668 bounding boxes of 1501 individuals captured using 6 cameras. These bounding boxes were obtained using the Deformable Part Model (DPM) ([Felzenszwalb et al., 2010](#)). The training set is made up of 751 identities with 12,936 images while the test set has 750 identities distinct from the one in the training set divided into query and gallery images.

MSMT17 ([Wei et al., 2018](#)): This is a very recent dataset which was carried out over a long period of time. This benchmark contains a total of 126,441 bounding boxes of 4101 identities captured using 15 cameras. The images vary in terms of location (outdoors, indoors), weather conditions (over a month), as well as different times of day (morning, noon, afternoon). The bounding boxes were obtained using Faster RCNN and corrected using labelers. Containing many variations makes this dataset challenging as well as a good benchmark to use.

DukeMTMC-reID ([Ristani et al., 2016](#)): The DukeMTMC-reID dataset is a small part of the bigger DukeMTMC dataset that is usually used for multi-target multi-camera tracking. It is taken from 8 different cameras, and the person bounding box is manually labeled. It is made up of 1404 different identities with 702 identities used for training and 702 other identities used for testing.

VERI-Wild ([Lou et al., 2019b](#)): The VERI-Wild dataset is a recently released large-scale vehicle re-identification dataset with 416,314 vehicle images of 40,671 IDs captured by 174 cameras. Vehicle recording is unconstrained and thus contains vehicles from different angles. This makes this dataset very challenging. The test set is divided into three subsets: small, medium, and large. We show our results on the small and medium subset since the large subset requires more memory.

EPFL Roundabout: In order to validate the use of our methods on images recorded from drone view, we build the EPFL Roundabout dataset by recording vehicle flow of five different locations using a drone. The division of the dataset is similar to VERI-Wild. To encourage efficient training, we set the training and testing ratio as 1:3. As a result, we obtain 85,268 images with 3479 IDs. An image example is shown in [Fig. 4](#). The dataset will be made publicly available.

CityFlow-ReID: The CityFlow-ReID dataset is a vehicle re-identification dataset with 56,277 bounding boxes of 666 IDs captured by 40 cameras. The training set is made up of 36,935 images of 333 different vehicle IDs. The test set contains the other 333 IDs across 18,290 gallery images and 1052 query images.

Table 1

Hyperparameters for the different datasets and methods. LR: learning rate, β : pre-factors for loss constraint, α : pre-factor for cross-entropy.

Parameters	Market-1501	DukeMTMC	MSMT17	VERI-Wild	CityFlow-ReID
LR_{cross}, α	$5 \times 10^{-4}, 2$	$2 \times 10^{-4}, 1$	$3 \times 10^{-4}, 4$	$1 \times 10^{-3}, 6$	$4 \times 10^{-4}, 1$
LR_{LS}, α	$5 \times 10^{-4}, 2$	$5 \times 10^{-4}, 5$	$3 \times 10^{-4}, 5$	$1 \times 10^{-3}, 6$	$4 \times 10^{-4}, 6$
LR_{CP}, α	$6 \times 10^{-4}, 3$	$6 \times 10^{-4}, 3$	$4 \times 10^{-4}, 5$	$1 \times 10^{-3}, 6$	$4 \times 10^{-4}, 5$
LR_{VIB}, α	$4 \times 10^{-4}, 6$	$6 \times 10^{-4}, 3$	$5 \times 10^{-4}, 6$	$1 \times 10^{-3}, 6$	$4 \times 10^{-4}, 5$
β_{LS}	0.1	0.1	0.1	0.1	0.05
β_{CP}	0.085	0.085	0.085	0.6	0.05
β_{VIB}	0.01	0.01	0.01	0.01	0.01

6.1. Evaluation protocol

For evaluation, we use the cumulative matching characteristic (CMC) and Mean Average Precision (mAP). These two metrics are the most popular evaluation metrics since re-identification systems should be able to output all the correct matches (mAP) in addition to having high accuracy at different ranks (CMC). During testing, for every query, there is a list of gallery images ordered in increasing order according to their L2 distance from this query.

6.2. Implementation details

The model was pre-trained using ImageNet. We do not add any layer to ResNet-50 when training both using label smoothing and confidence penalty except for a fully connected layer that outputs the different labels. When training the VIB algorithm, a fully-connected layer was added before the classification layer to output the mean and standard deviation which describe the distribution of the latent representations. A latent variable is then sampled from the predicted distribution. For all methods and datasets, hyperparameter tuning was performed for ResNet-50 in order to get the best possible accuracy.

Data Augmentation. We follow methods of data augmentations that are commonly used in the field of person re-identification. Since Market1501 uses DPM to obtain the bounding boxes, the images are initially randomly cropped. For all datasets, the inputs are resized to 256×128 . Before providing them to the network, a random rectangle, with pixel values randomly assigned between [0, 255], is erased (Zhong et al., 2017b) from the images, and the resulting images are flipped horizontally with a probability of 0.5. This makes the network more robust to the orientation of the agents in the image as well as occlusion. Each image is then normalized and standardized using the mean and standard deviation provided when using a model pretrained on ImageNet. These transformations were applied only for the training set.

Hyperparameter Tuning. Since the hyperparameters (e.g., learning rate, β , and α) we are trying to optimize have multiplicative effects on the training procedure, the best method is to perform a log-space search. This is due to two reasons. The parameter is not too sensitive such that there may not be too much difference with 10 and 15 compared to 10 and 1000. The other reason is that using logarithmic scales allows us to search over a bigger space quickly.

Training Procedure. The samples used to form the training batch are randomly sampled from the datasets. It does not require any special sampling such as the PK Sampling required by triplet loss (Schroff et al., 2015), which randomly samples P identities and then randomly K images for each identity to form a batch. The mini-batch has a size of 32 images. The model is trained for 300 epochs using AMSGrad (Reddi et al., 2018) for all datasets with the learning rate decaying by 10 at epoch 20 and 40. In order to make sure that all models were trained with the best parameters, we perform hyperparameter tuning, as discussed previously. The different hyperparameters for the different datasets are shown in Table 1.

Evaluation Procedure. For testing, the features that are extracted just before the last classification layer are used for the ranking process. The features for the queries and galleries are extracted and then compared to rank the gallery images relative to each query image. This is done when label smoothing or confidence penalty is used. When using the VIB loss, the network has an additional fully connected layer that outputs the mean and standard deviation for every latent dimension and a reparameterization trick that depends on random Gaussian noise. For ranking, we use the mean produced by the model as features for each image since this represents the average of the distribution over which the input image is mapped to. This is also due to the fact that the standard deviations tend to 1. To the best of our knowledge, using a latent representation sampled from a Gaussian parametrized by the predicted mean and standard deviation has not been tackled before for the person ad vehicle re-id task.

7. Results

In order to show both qualitative and quantitative results, we split our results into three parts. In Sections 7.1 and 7.2, we compare our proposed methods to published baseline results and state-of-the-art methods respectively. In Section 7.3, we investigate the effect of the three methods on the ranking process of person and vehicle re-id. Although these methods were tested on ResNet-50, other re-id models can benefit from their positive effect on the performance especially when dealing with visually similar pedestrians.

Table 2
Comparison with published ResNet-50 results on the Market-1501 and DukeMTMC-reID dataset.

Model	Market1501		DukeMTMC	
	mAP	Rank1	mAP	Rank1
ResNet-50 (Liu et al., 2018)	47.78	73.90	44.99	65.22
ResNet-50 (Shen et al., 2018b)	59.8	81.4	55.5	75.3
ResNet-50 (Saquib Sarfraz et al., 2018)	59.8	82.6	50.3	71.5
ResNet-50 (Chang et al., 2018)	66.0	84.3	48.6	71.6
ResNet-50 (Huang et al., 2018)	66.95	84.42	57.34	75.60
ResNet-50 (Kalayeh et al., 2018)	66.32	85.10	54.77	73.70
Our ResNet-50	70.2	87.5	59.6	78.6

7.1. Properly trained baseline

We compare our baseline to previously reported results of ResNet-50 on the Market-1501 and DukeMTMC-reID datasets. The published results reported in [Table 2](#) correspond to pre-trained ResNet-50 that used the cross-entropy loss similar to our method. As can be observed in [Table 2](#), there is a clear difference between our result and the results reported in published papers as well as amongst the published results themselves. Our properly trained baseline, which consists of a ResNet-50 model trained using a normal cross-entropy loss, was able to outperform all previous baselines. This table represents one of the many pitfalls that occurs when training a model. This is shown by the fact that papers that make use of exactly the same baseline have different results. This is usually due to the hyperparameters chosen. Another pitfall is that to compare different baselines and losses, the same hyperparameters are set. This is somewhat unfair since different baselines and losses optimize different parameters and in different ways thus requiring distinct hyperparameters. This is why we employ different learning rates for different datasets and methods as shown in [Table 1](#). As a result, we were able to achieve, using the baseline, around $\sim 5\%$ and $\sim 3\%$ relative increase in mAP and Rank-1 respectively for Market-1501, compared to ResNet-50 ([Kalayeh et al., 2018](#)). On DukeMTMC, our baseline achieves an mAP and Rank-1 of 59.6% and 78.6%, respectively.

7.2. Comparison with state-of-the-art

We evaluate our proposed confidence-based methods against recently published papers in person and vehicle re-id. Each of our methods is evaluated on five datasets: Market1501, DukeMTMC-reID, MSMT17, VERI-Wild, and EPFL Roundabout. We are able to reach state-of-the-art performance without any human-specific design and added complexity thus showing the importance of penalizing the confidence of a network in the task of re-identification. We also do not make use of data augmentation during the evaluation stage like DuATM ([Si et al., 2018](#)).

Evaluation on Market1501: As shown in [Table 3](#), the models were able to reach state-of-the-art results. In order to better understand the importance of penalizing confidence compared to other methods, it is important to note some distinct differences. Confidence penalty was able to outperform HAP2S ([Yu et al., 2018](#)) which tried to deal with hard samples by giving them higher weights. Moreover, Mancs ([Wang et al., 2018b](#)), which shows good performance, makes use of three different losses, attention layers, as well as a special sampling scheme. To compare our results with methods that include gallery-to-gallery information during inference, such as Deep Group RW ([Shen et al., 2018b](#)) and SGGNN ([Shen et al., 2018a](#)), we apply re-ranking to our three methods. We were able to outperform these methods with a significant relative increase in mAP($\sim 9\%$). As a result, we got state-of-the-art performance without the added complexity of learning new layers and parameters while tackling the problem stated in [Section 3](#).

Evaluation on DukeMTMC-reID: Similar to the Market-1501 dataset, we achieved competitive results in all proposed methods with confidence penalty resulting in the best improvement ([Table 3](#)). In addition to that, using Saquib Sarfraz et al.'s (2018) recent re-ranking method (ECN), we were able to get better results than PSE ([Saquib Sarfraz et al., 2018](#)) in both mAP and Rank-1. It is important to note that SPReID augments the training data of both DukeMTMC-reID and Market1501 with 10 datasets resulting in a large number of training samples which would improve the performance of the network.

Evaluation on MSMT17: Since this is a bigger dataset with many variations, it proved to be a challenging benchmark ([Wei et al., 2018](#)). Nonetheless, we were able to show a notable improvement over previous methods as well as over our own baseline ([Table 4](#)). Similarly, confidence penalty performed the best by achieving 68.6% in Rank-1 and 39.3% in mAP. By applying re-ranking, both Rank-1 and mAP are further improved to 75.3% and 59.1% respectively.

Evaluation on VERI-Wild: This dataset contains a large amount of images and is divided into small, medium, and large subsets. We evaluate our model on the small and medium subset and are able to achieve state-of-the-art performance ([Table 5](#)). Compared to the person re-id datasets, using the VIB method achieves the best performance on VERI-Wild. This might be due to the fact that vehicles face the problems of similar appearance and different IDs more frequently. Thus, a more strict method of penalizing certainty is required compared to confidence penalty and label smoothing.

Table 3

Comparison with state-of-the-art methods on Market-1501 and DukeMTMC-reID. (!R): uses model different than ResNet, (R): uses ResNet-50, ECN: Expanded Cross Neighborhood Re-Ranking ([Saquib Sarfraz et al., 2018](#)), “RR”: k-reciprocal re-ranking ([Zhong et al., 2017a](#)), Xent: Softmax.

Model	Market1501		DukeMTMC	
	mAP	Rank1	mAP	Rank1
CamStyle (R) (Zhong et al., 2018)	71.55	89.49	57.61	78.32
HAP2S_E + Xent(R) (Yu et al., 2018)	74.49	89.73	62.62	79.08
DuATM(!R) (Si et al., 2018)	75.22	89.96	63.14	81.46
MLFN (!R) (Chang et al., 2018)	74.3	90.0	62.8	81.0
Shen et al.(R) (Shen et al., 2018b)	75.3	90.1	63.2	80.3
PSE(R) (Saquib Sarfraz et al., 2018) + ECN	84.0	90.3	79.8	85.2
DaRe(IR) (Wang et al., 2018a) + RR	86.7	90.9	80.0	84.4
SPReID ^{w//s} (IR) (Khalayeh et al., 2018) ^a	78.66	90.97	65.66	81.73
HA-CNN (IR) (Li et al., 2018)	75.7	91.2	63.8	80.5
DuATM(IR) (Si et al., 2018) ^b	76.62	91.42	64.58	81.82
SPReID ^{comb} (IR) (Khalayeh et al., 2018) ^a	79.67	91.45	68.78	83.3
P-Aligned (!R) (Suh et al., 2018)	79.6	91.7	69.3	84.4
SGGNN(R) (Shen et al., 2018a)	82.8	92.3	68.2	81.1
Deep Group RW(R) (Shen et al., 2018b)	82.5	92.7	66.4	80.7
Mancs(R) (Wang et al., 2018b)	82.3	93.1	71.8	84.9
DNN + CRF(R) (Chen et al., 2018)	81.6	93.5	69.5	84.9
PCB + RPP (Sun et al., 2018)	81.6	93.8	69.2	83.3
P-Aligned (!R) (Suh et al., 2018) + RR	89.9	93.4	83.9	88.3
Our ResNet	70.7	87.2	59.6	78.6
Our ResNet(VIB)	76.1	90.2	62.4	80.7
Our ResNet(LS)	76.7	91.0	64.4	82.7
Our ResNet(CP)	78.2	91.4	66.8	83.9
Our ResNet + RR	85.7	89.7	78.5	83.4
Our ResNet(VIB) + RR	88.6	91.8	79.0	84.3
Our ResNet(LS) + RR	89.1	92.2	82.2	86.6
Our ResNet(CP) + RR	90.0	92.6	83.5	87.4
Our ResNet(VIB) + ECN	88.2	92.0	78.9	85.1
Our ResNet(LS) + ECN	89.4	92.7	83.2	86.9
Our ResNet(CP) + ECN	90.1	93.1	84.1	88.5

^aUses combination of 10 datasets for training.

^bUses data augmentation during evaluation stage.

Table 4

Comparison with state-of-the-art on the MSMT17 dataset.

MSMT17	mAP	Rank1	Rank10
GoogleNet (Wei et al., 2018)	23.0	47.6	71.8
PDC (Wei et al., 2018)	29.7	58.0	79.4
GLAD (Wei et al., 2018)	34.0	61.4	81.6
Our ResNet	31.8	59.3	80.2
Our ResNet(VIB)	35.1	66.2	84.1
Our ResNet(LS)	36.9	66.8	84.9
Our ResNet(CP)	39.3	68.6	85.3
Our ResNet + RR	49.8	65.7	79.8
Our ResNet(VIB) + RR	55.4	73.3	84.7
Our ResNet(LS) + RR	57.1	73.7	85.3
Our ResNet(CP) + RR	59.1	75.3	85.8

Evaluation on EPFL Roundabout: Since the images are recorded from a drone view, vehicles are small and lack certain details. This makes it more challenging to discriminate between different vehicles and a result suffers more from the problem discussed in Section 3. However, we were able to drastically improve the performance of the baseline by making use of confidence penalty, label smoothing, and VIB. Confidence Penalty, in this case shows, the biggest relative improvement of around 35.2% and 13.6% in both mAP and Rank1 respectively.(Table 6).

Evaluation on CityFlow-ReID: In order to properly train our models on this dataset, the images are resized to 512×256 . Since the dataset does not have any validation set, we only use 85% of the training images while the remaining 15% are used to determine when to stop training to prevent overfitting. As shown in Table 7, our methods outperform previous models and improve the baseline results despite being trained on part of the training set. We also report the performance of our baseline since we obtain a $\sim 7.3\%$ relative increase in rank1 compared to the ResNet baseline reported in the cityflow paper ([Tang et al., 2019](#)). VIB leads to the best

Table 5
Comparison with state-of-the-art on VERI-Wild dataset.

Model	Small		Medium	
	mAP	Rank1	mAP	Rank1
GSTE (Bai et al., 2018)	31.4	60.5	26.18	52.12
VERI-Wild (Lou et al., 2019b)	35.1	64.0	29.8	57.8
Our ResNet	45.7	82.4	41.3	78.4
Our ResNet(LS)	57.7	84.6	57.2	85.5
Our ResNet(CP)	67.5	90.2	61.8	87.0
Our ResNet(VIB)	74.1	92.1	68.5	89.7

Table 6
Improved performance of our methods on EPFL Roundabout ReID.

Model	EPFL Roundabout ReID	
	mAP	Rank1
ResNet	41.5	75.0
Our ResNet(LS)	55.9	84.4
Our ResNet(CP)	56.1	85.2
Our ResNet(VIB)	52.7	82.5

Table 7
Improved performance of our methods on CityFlow-ReID test-set.

Model	CityFlow-ReID	
	mAP	Rank1
ResNet (Tang et al., 2019)	25.5	41.3
ResNeXt101 (Tang et al., 2019)	26.6	42.4
SEResNeXt50 (Tang et al., 2019)	26.8	45.2
Our ResNet	22.42	44.3
Our ResNet(LS)	25.84	51.81
Our ResNet(CP)	26.92	54.28
Our ResNet(VIB)	28.83	53.90

performance with $\sim 28.6\%$ and 21.7% relative increase in mAP and Rank-1, respectively, compared to our baseline. All our methods result in an increase in mAP and lead to a drastic increase in Rank-1 compared to previous methods, emphasizing the improved features extracted by our models.

7.3. Effect of proposed methods

In addition to achieving state-of-the-art performance, it is also important to understand the effect of these three methods on the ranking process. All three methods aim at allowing the network to share some representation among different classes. This prevents the network from focusing on undesirable information when separating very similar-looking pedestrians. To show this effect, we compare the confidence penalty model against the baseline model since it resulted in the best performance in person re-id. As can be seen, the test samples presented in Figs. 5 and 6 are difficult to rank even for an observer. This confirms the intrinsic difficulty of person and vehicle re-id stated in Section 3. When confidence penalty is not used for training, the network focuses on unimportant variations between the images. For instance, in both sets of samples (Fig. 5), the incorrect gallery images are very similar to the query image despite belonging to a different person. The baseline links the query image to the gallery images by possibly focusing on the background, shirt color, posture, and body rotation of the pedestrian in question. The same applies for vehicles. These characteristics are typically features that can confuse the model leading to wrong identification. Adding the confidence penalty is observed to remedy this challenge, as can be seen for all test samples provided. Adding the confidence penalty helps the model capture the subtle differences between multiple pedestrians that the baseline tends to misidentify. These are ideal examples of why confidence penalty drastically improved re-identification compared to less significant improvements in other computer vision tasks.

8. Ablation study

8.1. Combination of penalties

Since the addition of the three methods to the baseline leads to an improvement in performance, one might wonder the effect of the different combination of these methods. As shown in Table 8, applying only confidence penalty (CP) leads to the best result.



Fig. 5. Qualitative comparison of using confidence penalty on unseen Market1501 test samples. The gallery images are ranked according to L2 distance (top-5, left to right). Red frame indicates wrong ID while green frame indicates correct ID compared to the query. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

Table 8
Ablation study of different penalty combinations.

Model	Market1501		DukeMTMC	
	mAP	Rank1	mAP	Rank1
ResNet	70.7	87.2	59.6	78.6
Our ResNet (CP)	78.2	91.4	66.8	83.9
Our ResNet (VIB + LS)	76.8	90.4	64.1	80.9
Our ResNet (VIB + CP)	76.7	90.2	63.6	80.7
Our ResNet (LS + CP)	77.2	91.3	63.9	82.0
Our ResNet (VIB + LS + CP)	76.1	89.6	63.8	81.5

Intuitively, the combinations of the methods increases the restriction on the model preventing it from learning useful representations. Meanwhile, adding these methods together still leads an improved performance compared to the baseline which again affirms the beneficial effect that they have. The hyperparameters used to train these models are reported in [Table 9](#). They were determined using the technique explained in Section 6.2 where log-space search was performed over the hyperparameters.

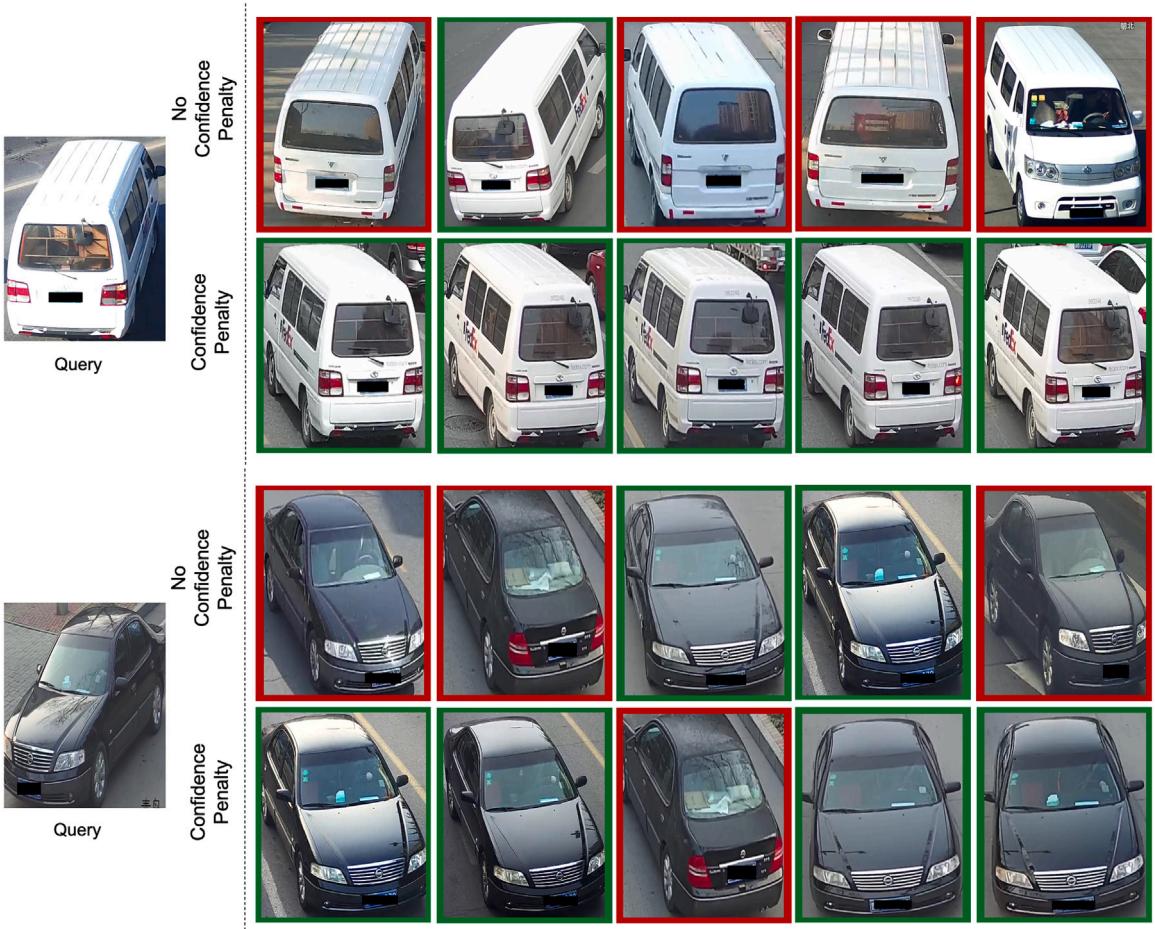


Fig. 6. Qualitative comparison of using confidence penalty on unseen VERI-Wild test samples. The gallery images are ranked according to L2 distance (top-5, left to right). Red frame indicates wrong ID while green frame indicates correct ID compared to the query. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

Table 9
Hyperparameters of different penalty combinations.

Model	Market1501			DukeMTMC		
	LR	β	α	LR	β	α
Our ResNet (VIB + LS)	1×10^{-3}	$\beta_{VIB} = 0.01$ $\beta_{LS} = 0.1$	4	6×10^{-4}	$\beta_{VIB} = 0.01$ $\beta_{LS} = 0.1$	6
Our ResNet (VIB + CP)	3×10^{-4}	$\beta_{VIB} = 0.01$ $\beta_{CP} = 0.05$	4	4×10^{-4}	$\beta_{VIB} = 0.01$ $\beta_{CP} = 0.1$	5
Our ResNet (LS + CP)	8×10^{-4}	$\beta_{LS} = 0.1$ $\beta_{CP} = 0.2$	4	2×10^{-4}	$\beta_{LS} = 0.1$ $\beta_{CP} = 0.1$	4
Our ResNet (VIB + LS + CP)	8×10^{-4}	$\beta_{VIB} = 0.01$ $\beta_{LS} = 0.1$ $\beta_{CP} = 0.06$	4	6×10^{-4}	$\beta_{VIB} = 0.01$ $\beta_{LS} = 0.1$ $\beta_{CP} = 0.1$	5

8.2. Alpha-beta parameters

The alpha and beta in each method, in addition to the learning rate and other hyperparameters, have an impact on the performance of the re-identification model. In order to better understand the effect and obtain an approximate of the best intervals for the hyperparameters, we performed a grid search over different alphas and betas on Market1501, DukeMTMC-reID, and Cityflow-reID. To also show that different learning rates share similar alphas and betas, each of the datasets are trained using a different learning rate while varying over the same alphas and betas. As shown in Fig. 7, each method, irrespective of the dataset, obtains the best performance over a specific alpha and beta interval. It is important to note that none of the graphs begin at a beta equal

Table 10

The performance and complexity of the different methods. Test time (ms) measures the time required to obtain the features of a single image. GMACs measures the number of multiply-accumulates per second.

Model	Test time (ms)	GMACs
ResNet	5.68	2.683
Our ResNet (LS)	5.68	2.683
Our ResNet (CP)	5.68	2.683
Our ResNet (VIB)	5.85	2.688

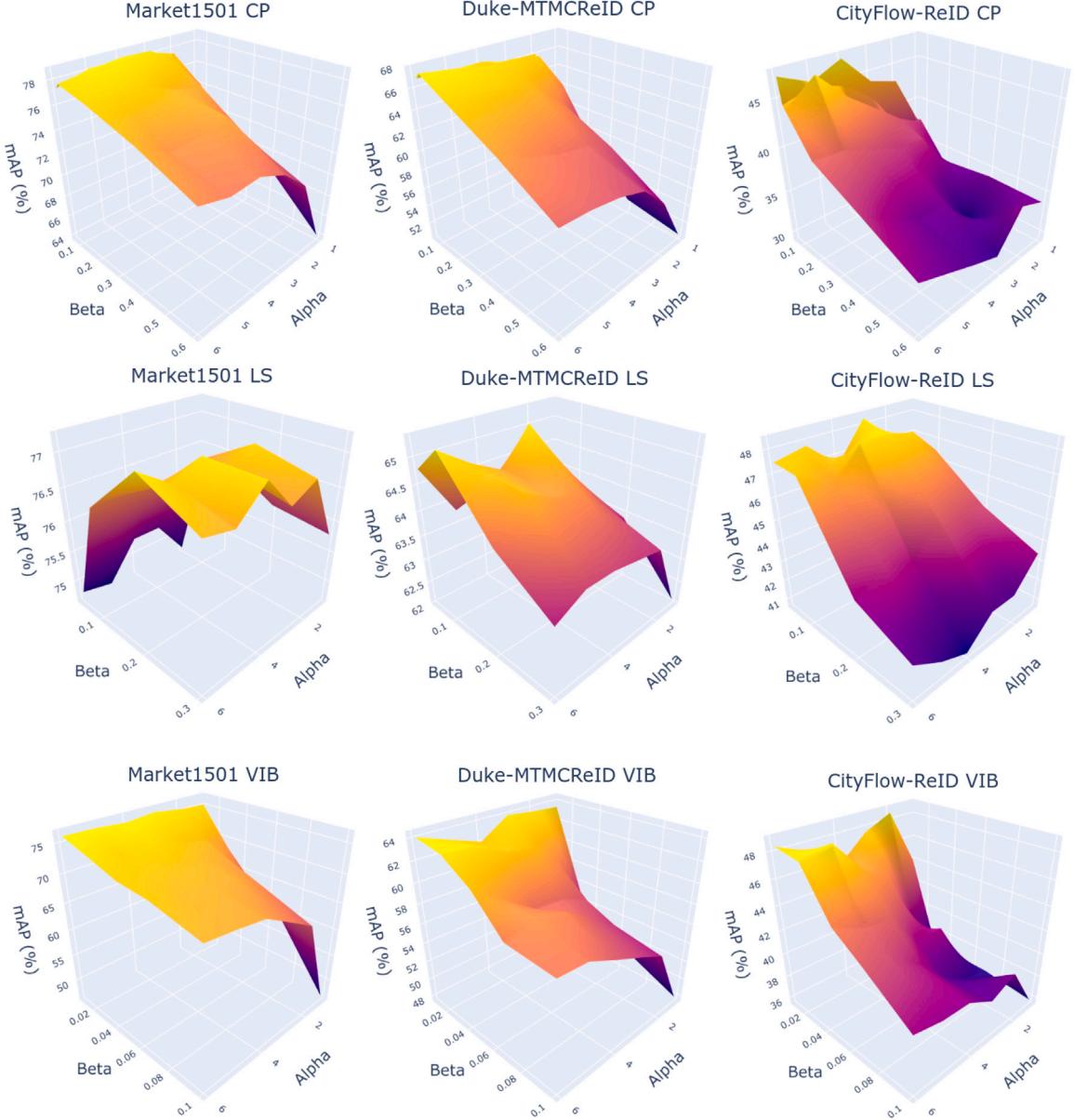


Fig. 7. Variation of the mAP (%) versus alpha and beta for the three different methods (CP, LS, VIB) on Market1501, Duke-MTMCReID, and CityFlow-ReID.

to zero. For instance, confidence penalty has shown high mAP at beta and alpha equal to 0.1 and 5 respectively. Label smoothing resulted in the best performance when setting beta equal to 0.1 or 0.2 and alpha to 3 or 4 while VIB achieved the highest mAP at beta and alpha equal to 0.01 and 6, respectively. It is important to notice that different datasets have different z ranges and that different alpha and beta combinations can lead to similar performance.



Fig. 8. A failure case from Market1501 dataset. The query and gallery images do not contain any distinctive feature to discern them from each other. After thorough investigation, we believe these images are mislabeled and should pertain to the same person since the person is wearing the same clothe and shoes as well as carries the same bag on the side.

Table 11

Effect of CP and LS on PCB, HA-CNN and BFE. Xent: Cross-Entropy, G: using global features, L: using local features.

Model	Market1501		DukeMTMC	
	mAP	Rank1	mAP	Rank1
PCB (own) (Sun et al., 2018)	78.5	92.4	69.6	83.8
Our PCB(LS)	77.1	91.7	69.0	84.6
Our PCB(CP)	78.8	93.2	70.1	84.3
PCB + RPP (own) (Sun et al., 2018)	81.5	93.3	71.3	84.4
Our PCB + RPP(LS)	80.9	92.8	71.3	85.3
Our PCB + RPP(CP)	81.6	93.3	71.5	84.2
HA-CNN (Li et al., 2018)	75.7	91.2	63.8	80.5
HA-CNN(own)	78.6	91.6	67.3	83.1
Our HA-CNN(LS)	80.1	92.0	68.4	83.6
Our HA-CNN(CP)	81.1	92.5	69.7	83.8
HA-CNN(G)	72.3	87.9	60.0	78.7
Our HA-CNN(CP + G)	77.3	91.5	63.4	82.0
HA-CNN(L)	73.8	89.5	62.3	81.0
Our HA-CNN(CP + L)	75.5	91.2	62.8	81.0
BFE (Xent) (Dai et al., 2018)	83.7	93.5	73.5	86.4
Our BFE (CP)	85.7	94.2	76.1	88.6

Table 12
Effect of confidence penalty on NuScenes-ReID.

Model	nuScenes-ReID	
	mAP	Rank1
ResNet	61.7	66.3
Our ResNet(LS)	66.0	70.4
Our ResNet(CP)	70.7	74.4
Our ResNet(VIB)	67.7	70.4

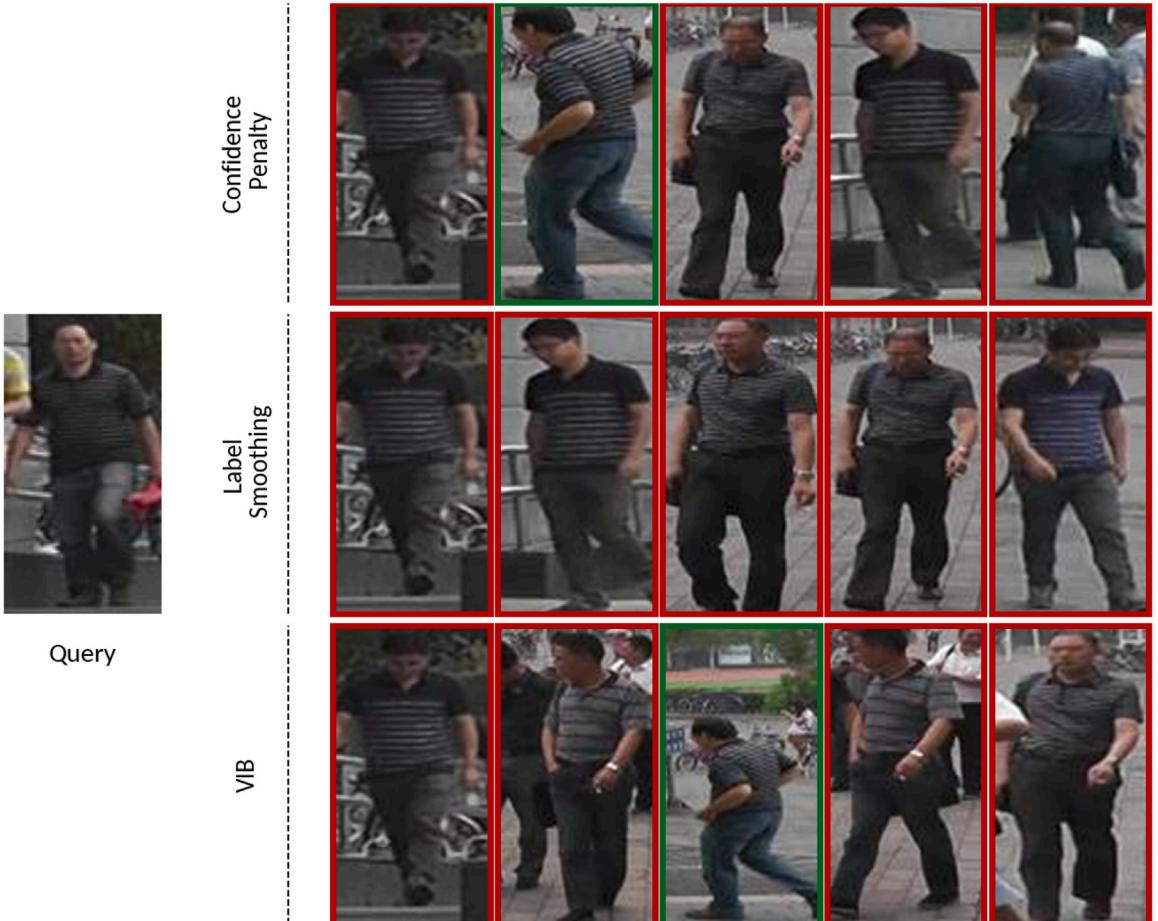


Fig. 9. A failure case from Market1501 dataset. Since we are using, ResNet50, that performs a global average pooling at the last layer as well as downscale the input image, certain details can be diluted with other more visible features. For instance, all the shirts in the above image are striped but the network fails to distinguish the frequency and size of these stripes. A human observer might still be able to notice such differences.

8.3. Computational complexity

During testing, all methods do not make use of the final fully connected layer responsible to predict the ID of every input image. Moreover, the methods use the same input size and same base network, Resnet50, with the exception of VIB requiring an additional fully connected layer (check Fig. 3). The time required to extract the features from a single image as well as the computational complexity of each method are shown in Table 10. Since most modern hardwares use FMA instructions which perform fused multiply-add (FMA) operations, we report the computational complexity in GMACs (giga multiply-accumulates per second). Roughly, GMACs is equal to half a GFLOPs. Since confidence penalty and label smoothing only affect the loss function during training, they do not increase the complexity of the model during testing. The additional fully connected layer in VIB has negligible effect on the complexity as well.

8.4. Failure cases

Confidence penalty, label smoothing, and VIB tackle a problem in re-identification models discussed in Section 3. We still experience certain failure cases due to the different difficulties faced by a visual re-identification model, such as misalignment



Fig. 10. A failure case from VERI-Wild dataset. The images above do not contain any detail on the vehicles to distinguish them from each other since they are of the same model and same color. A human annotator would also fail to correctly rank the images.

as well as dataset quality. Augmenting these methods to other state-of-the-art approaches that deal with such challenges will lead to an improved and more complete model, as shown in Section 8.5. We show two failure cases from each of Market1501 (Figs. 8–9) and VERI-Wild (Figs. 10–11) datasets and provide an analysis behind these failures. The Market1501 dataset appears to include different labels for the same person which penalizes our method when correctly ranking the gallery based on the query image (Fig. 8). Several vehicles on the roads are of the same model and color and do not have distinctive features. As a result, our models are unable to discern between such examples in VERI-Wild (Fig. 10), especially since license plates are hidden. Another failure might be due to the fact that we are downsampling the image using ResNet50 and then applying a global average pooling to obtain the features. This can cause certain features to dilute with other more visible features (Figs. 9 and 11). Making use of more complicated methods can help deal with this problem while still benefitting from the improvements our methods provide, as shown in Table 11.

8.5. Penalties on state-of-the-art methods

To study the effect of our best method on the different state-of-the-art methods, we test confidence penalty on PCB (Sun et al., 2018), HACNN (Li et al., 2018), and BFE (Dai et al., 2018). It is important to note that PCB divides the representation into multiple features that are then used for identification. These features are referred to as local features since they are spatially local to a certain region of the input. In addition to local features, HACNN also extracts a global representation from the whole input. In comparison, BFE performs person re-identification using only a global representation. Table 11 shows that the relative gain in performance for PCB is between +0.4 to +1%. One intuition behind the limited gain is that PCB divides its features into multiple parts (local features) before applying global average pooling. This allows the representations to be fine-grained focusing on the details that differentiate visually similar inputs. To analyze this, we test CP on HA-CNN that uses global and local features. Its performance is improved ($\sim +3\%$ mAP) to exceed PCB and to have on-par results with PCB+RPP. Also, CP has a bigger effect on global features than on local features in HA-CNN (Table 11) confirming our reasoning. The performance of another method, BFE, that uses only global features is also improved using CP compared to cross-entropy (Xent).

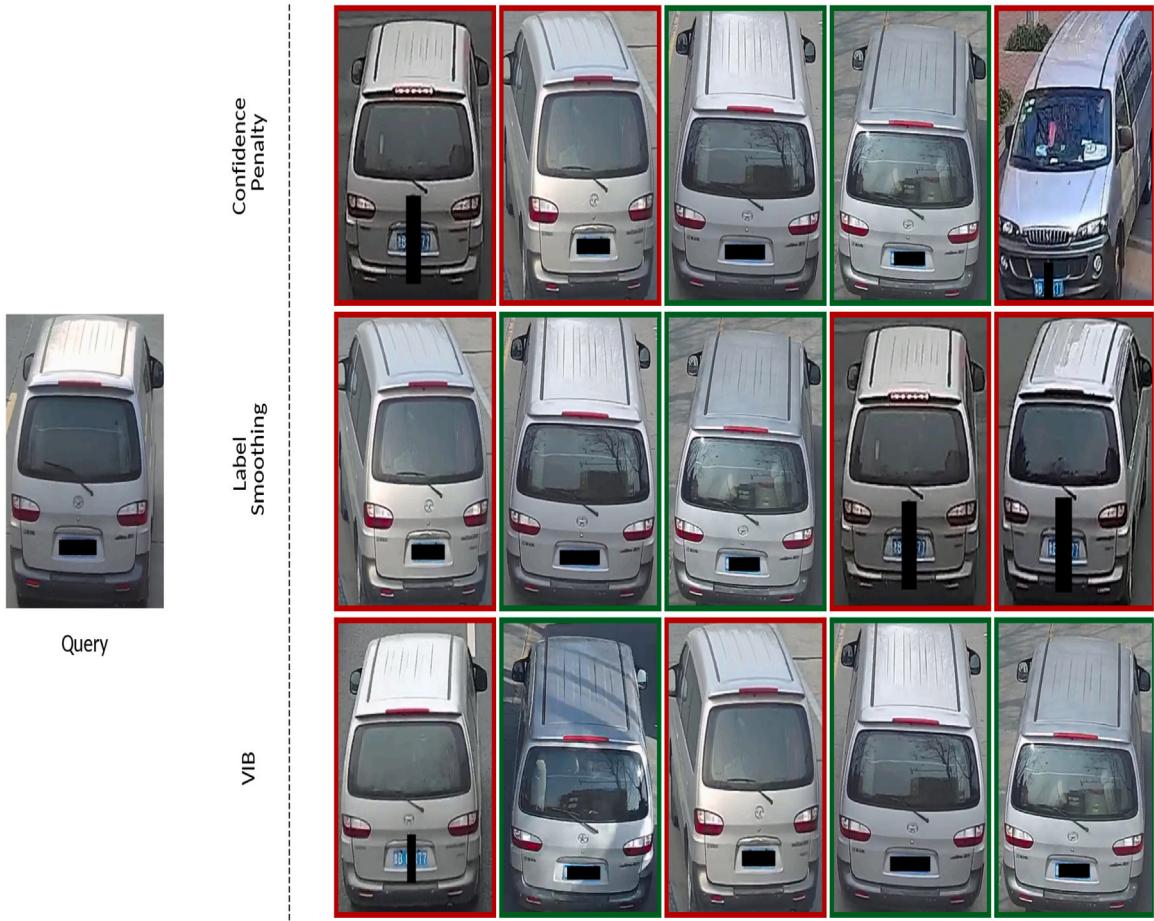


Fig. 11. A failure case from VERI-Wild dataset. Since we are using, ResNet50, that performs a global average pooling at the last layer as well as downscale the input image, certain details can be diluted with other more visible features. For instance, the two writings on the right side below the lamp might be considered as one when passing the image through the network. A human annotator would have noticed these details.

8.6. Vehicle ego-centric pedestrian re-identification

The dataset that were used to evaluate our methods are collected from cameras mounted in specific locations or using drones. In order to test the effectiveness of our method in a scenario where the camera is mounted on a car, we make use of the nuScenes dataset. This is the first dataset to contain all sensor information that a full autonomous vehicles requires from RGB camera, radars to lidars. Since the pedestrians are tracked across images and different cameras, we can build a re-identification dataset similar to Market1501. We call the resulting dataset, nuScenes-ReID. This is a challenging dataset since images are recorded from the view of a moving car resulting in different pedestrian sizes. As shown in Table 12, applying confidence penalty to the baseline significantly improves the performance (from 59.1%mAP to 70.9%). The code to build this re-id dataset will be made public.

9. Conclusions

An important task of transportation research is better analyzing and understanding traffic flow. Visual re-identification, the task of association similar agents, can aid in this goal. Thus, in this work we aim to deal with certain challenges that plague this task. First, we emphasize an intrinsic characteristic of person and vehicle re-identification that poses a problem to the network being trained. The classes that these re-id task try to separate are not as easy as separating cats and dogs. Different agents with different identities can have very similar appearances. We have demonstrated that three methods, that reduce a model's confidence, are able to deal with this problem while achieving state-of-the-art results. Confidence penalty proved to be the best and most lightweight amongst the three different methods. In addition, it is interesting to note that VIB is able to achieve similar results while using smaller representations. Both label smoothing and confidence penalty use a representation of dimension 2048 while VIB uses a representation of dimension 1024. These three methods can be leveraged to improve the performance of previous re-id methods as well. It remains an exciting future work to study their effect on other image retrieval and clustering tasks.

With the ability to identify pedestrians and cars across both time and space, this allows us to better understand how they move from one place to the other without going through the manual and expensive way of using surveys or interviews. Some methods have also been developed to estimate the OD matrices from existing observed traffic flows. These methods, however, require the collected data to be large and representative of the real distribution. This is where re-identification can play a major role. CCTV cameras already placed around entry and exit of different transportation stops can be used to associate agents that pass through them. For instance, a person can be detected entering a specific train and then this detection can be associated to the same person exiting at a different station and at different times. This task provides an automatic method of collecting data about passenger movements and thus allowing us to build an accurate OD matrix that can be used for different transportation tasks.

CRediT authorship contribution statement

George Adaimi: Conceptualization, Methodology, Software, Investigation, Validation, Visualization, Writing - original draft. **Sven Kreiss:** Supervision, Writing - review & editing. **Alexandre Alahi:** Supervision, Writing - review & editing, Resources, Project administration.

Acknowledgment

This work was supported by the Swiss National Science Foundation Spark fund (190677) and EPFL Interdisciplinary fund. We also thank our lab members and reviewers for their valuable comments.

Appendix A. Supplementary data

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.trc.2021.103067>. It contains qualitative comparisons of confidence penalty, label smoothing, and VIB on Market-1501 and VERI-Wild.

References

- Abdulhai, Baher, Tabib, Seyed M., 2003. Spatio-temporal inductance-pattern recognition for vehicle re-identification. *Transp. Res. C* 11 (3–4), 223–239.
- Ahmed, Ejaz, Jones, Michael, Marks, Tim K., 2015. An improved deep learning architecture for person re-identification. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- Alahi, Alexandre, Bierlaire, Michel, Kunt, Murat, 2008. Object detection and matching with mobile cameras collaborating with fixed cameras. In: ECCV Workshop on Multi-Camera and Multi-Modal Sensor Fusion Algorithms and Applications.
- Alahi, Alexandre, Bierlaire, Michel, Vandergheynst, Pierre, 2014. Robust real-time pedestrians detection in urban environments with low-resolution cameras. *Transp. Res. C* 39, 113–128.
- Alahi, Alexandre, Ramanathan, Vignesh, Fei-Fei, Li, 2017a. Tracking millions of humans in crowded spaces. Book Chapter on Group and Crowd Behavior. Academic Press, pp. 115–135.
- Alahi, Alexandre, Wilson, Judson, Fei-Fei, Li, Savarese, Silvio, 2017b. Unsupervised camera localization in crowded spaces. In: IEEE International Conference on Robotics and Automation (ICRA).
- Alemi, Alexander A., Fischer, Ian, Dillon, Joshua V., Murphy, Kevin, 2016. Deep variational information bottleneck. CoRR, [abs/1612.00410](https://arxiv.org/abs/abs/1612.00410).
- Bai, Yan, Lou, Yihang, Gao, Feng, Wang, Shiqi, Wu, Yuwei, Duan, Ling-Yu, 2018. Group-sensitive triplet embedding for vehicle reidentification. *IEEE Trans. Multimed.* 20 (9), 2385–2399.
- Barmounakis, Emmanouil, Geroliminis, Nikolas, 2020. On the new era of urban traffic monitoring with massive drone data: The pNEUMA large-scale field experiment. *Transp. Res. C* 111, 50–71.
- Chang, Xiaobin, Hospedales, Timothy M., Xiang, Tao, 2018. Multi-level factorisation net for person re-identification. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- Chen, Weihua, Chen, Xiaotang, Zhang, Jianguo, Huang, Kaiqi, 2017. Beyond triplet loss: A deep quadruplet network for person re-identification. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- Chen, Dapeng, Xu, Dan, Li, Hongsheng, Sebe, Nicu, Wang, Xiaogang, 2018. Group consistent similarity learning via deep crf for person re-identification. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- Cheng, De, Gong, Yihong, Zhou, Sanping, Wang, Jinjun, Zheng, Nanning, 2016. Person re-identification by multi-channel parts-based cnn with improved triplet loss function. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- Corvee, Etienne, Bremond, Francois, Thonnat, Monique, et al., 2010a. Person re-identification using haar-based and dcd-based signature. In: 2010 7th IEEE International Conference on Advanced Video and Signal Based Surveillance. IEEE, pp. 1–8.
- Corvee, Etienne, Bremond, Francois, Thonnat, Monique, et al., 2010b. Person re-identification using spatial covariance regions of human body parts. In: 2010 7th IEEE International Conference on Advanced Video and Signal Based Surveillance. IEEE, pp. 435–440.
- Dai, Zuozhuo, Chen, Mingqiang, Zhu, Siyu, Tan, Ping, 2018. Batch feature erasing for person re-identification and beyond. ArXiv, [abs/1811.07130](https://arxiv.org/abs/abs/1811.07130).
- Felzenszwalb, P.F., Girshick, R.B., McAllester, D., Ramanan, D., 2010. Object detection with discriminatively trained part-based models. *IEEE Trans. Pattern Anal. Mach. Intell.* 32 (9), 1627–1645.
- Gaddigoudar, P.K., Balihalli, T.R., Ijantkar, S.S., Iyer, N.C., Maralappanavar, S., 2017. Pedestrian detection and tracking using particle filtering. In: 2017 International Conference on Computing, Communication and Automation (ICCA). pp. 110–115.
- García, J., Martinel, N., Micheloni, C., Gardel, A., 2015. Person re-identification ranking optimisation by discriminant context information analysis. In: 2015 IEEE International Conference on Computer Vision (ICCV). pp. 1305–1313.
- Gheissari, N., Sebastian, T.B., Hartley, R., 2006. Person reidentification using spatiotemporal appearance. In: 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06), Vol. 2. pp. 1528–1535.
- Guo, Lie, Li, Linhui, Zhao, Yibing, Zhao, Zongyan, 2016. Pedestrian tracking based on camshift with Kalman prediction for autonomous vehicles. *Int. J. Adv. Robot. Syst.* 13 (3), 120.
- Hadsell, R., Chopra, S., LeCun, Y., 2006. Dimensionality reduction by learning an invariant mapping. In: 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06), Vol. 2. pp. 1735–1742.

- Held, David, Levinson, Jesse, Thrun, Sebastian, Savarese, Silvio, 2014. Combining 3D shape, color, and motion for robust anytime tracking.. In: *Robotics: Science and Systems*.
- Huang, Houjing, Li, Dangwei, Zhang, Zhang, Chen, Xiaotang, Huang, Kaiqi, 2018. Adversarially occluded samples for person re-identification. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Huang, Timothy, Russell, Stuart, 1997. Object identification in a Bayesian context. In: *Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence - Volume 2*. In: IJCAI'97, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, pp. 1276–1282.
- Hue, Carine, Le Cadre, J.-P., Pérez, Patrick, 2002. Tracking multiple objects with particle filtering. *IEEE Trans. Aerosp. Electron. Syst.* 38 (3), 791–812.
- Kalayeh, Mahdi M., Basaran, Emrah, Gökmén, Muhittin, Kamasak, Mustafa E., Shah, Mubarak, 2018. Human semantic parsing for person re-identification. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Karianakis, Nikolaos, Liu, Zicheng, Chen, Yinpeng, Soatto, Stefano, 2017. Person depth reid: Robust person re-identification with commodity depth sensors. *CoRR*, abs/1705.09882.
- Kingma, Diederik P., Welling, Max, 2014. Auto-encoding variational Bayes. *CoRR*, abs/1312.6114.
- Krishnakumari, Panchamy, van Lint, Hans, Djukic, Tamara, Cats, Oded, 2019. A data driven method for OD matrix estimation. *Transp. Res. C*.
- Kullback, Solomon, Leibler, Richard A., 1951. On information and sufficiency. *Ann. Math. Stat.* 22 (1), 79–86.
- Leng, Qingming, Hu, Ruimin, Liang, Chao, Wang, Yimin, Chen, Jun, 2015. Person re-identification with content and context re-ranking. *Multimedia Tools Appl.* 74 (17), 6989–7014.
- Li, Dangwei, Chen, Xiaotang, Zhang, Zhang, Huang, Kaiqi, 2017a. Learning deep context-aware features over body and latent parts for person re-identification. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Li, Yuqi, Li, Yanghao, Yan, Hongfei, Liu, Jiaying, 2017b. Deep joint discriminative learning for vehicle re-identification and retrieval. In: *2017 IEEE International Conference on Image Processing (ICIP)*. IEEE, pp. 395–399.
- Li, Hui, Liu, Yun, Wang, Chuanxu, Zhang, Shujun, Cui, Xuehong, 2016. Tracking algorithm of multiple pedestrians based on particle filters in video sequences. *Comput. Intell. Neurosci.* 2016.
- Li, Wei, Zhao, Rui, Xiao, Tong, Wang, Xiaogang, 2014. Deepreid: Deep filter pairing neural network for person re-identification. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Li, Wei, Zhu, Xiatian, Gong, Shaogang, 2018. Harmonious attention network for person re-identification. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Liu, Xinchen, Liu, Wu, Ma, Huadong, Fu, Huiyuan, 2016a. Large-scale vehicle re-identification in urban surveillance videos. In: *2016 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, pp. 1–6.
- Liu, Xinchen, Liu, Wu, Mei, Tao, Ma, Huadong, 2016b. A deep learning-based approach to progressive vehicle re-identification for urban surveillance. In: *European Conference on Computer Vision*. Springer, pp. 869–884.
- Liu, Jinxian, Ni, Bingbing, Yan, Yichao, Zhou, Peng, Cheng, Shuo, Hu, Jianguo, 2018. Pose transferrable person re-identification. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Liu, Hongye, Tian, Yonghong, Yang, Yaowei, Pang, Lu, Huang, Tiejun, 2016. Deep relative distance learning: Tell the difference between similar vehicles. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 2167–2175.
- Lou, Yihang, Bai, Yan, Liu, Jun, Wang, Shiqi, Duan, Ling-Yu, 2019a. Embedding adversarial learning for vehicle re-identification. *IEEE Trans. Image Process.* 28 (8), 3794–3807.
- Lou, Yihang, Bai, Yan, Liu, Jun, Wang, Shiqi, Duan, Lingyu, 2019b. Veri-wild: A large dataset and a new method for vehicle re-identification in the wild. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Mazzon, Riccardo, Tahir, Syed Fahad, Cavallaro, Andrea, 2012. Person re-identification in crowd. *Pattern Recognit. Lett.* 33 (14), 1828–1837.
- Pereyra, Gabriel, Tucker, George, Chorowski, Jan, Kaiser, Lukasz, Hinton, Geoffrey E., 2017. Regularizing neural networks by penalizing confident output distributions.
- Rama Varior, Rahul, Shuai, Bing, Lu, Jiwen, Xu, Dezhi, Wang, Gang, 2016. A Siamese Long Short-Term Memory Architecture for Human Re-identification, Vol. 9911. pp. 135–153.
- Reddi, Sashank J., Kale, Satyen, Kumar, Sanjiv, 2018. On the convergence of adam and beyond. In: *International Conference on Learning Representations*.
- Ristani, Ergys, Solera, Francesco, Zou, Roger, Cucchiara, Rita, Tomasi, Carlo, 2016. Performance measures and a data set for multi-target, multi-camera tracking. In: *European Conference on Computer Vision*. Springer, pp. 17–35.
- Saqib Sarfraz, M., Schumann, Arne, Eberle, Andreas, Stieffelhagen, Rainer, 2018. A pose-sensitive embedding for person re-identification with expanded cross neighborhood re-ranking. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Schneider, Benjamin, 2018. The High Cost of Global Traffic Jams. CityLab.
- Schroff, Florian, Kalenichenko, Dmitry, Philbin, James, 2015. Facenet: A unified embedding for face recognition and clustering. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Shen, Yantao, Li, Hongsheng, Yi, Shuai, Chen, Dapeng, Wang, Xiaogang, 2018a. Person re-identification with deep similarity-guided graph neural network. In: *The European Conference on Computer Vision (ECCV)*.
- Shen, Yantao, Xiao, Tong, Li, Hongsheng, Yi, Shuai, Wang, Xiaogang, 2018b. End-to-end deep kronecker-product matching for person re-identification. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Si, Jianlou, Zhang, Honggang, Li, Chun-Guang, Kuen, Jason, Kong, Xiangfei, Kot, Alex C., Wang, Gang, 2018. Dual attention matching network for context-aware feature sequence based person re-identification. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Song, Chunfeng, Huang, Yan, Ouyang, Wanli, Wang, Liang, 2018. Mask-guided contrastive attention model for person re-identification. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Suh, Yumin, Wang, Jingdong, Tang, Siyu, Mei, Tao, Mu Lee, Kyoung, 2018. Part-aligned bilinear representations for person re-identification. In: *The European Conference on Computer Vision (ECCV)*.
- Sun, Yifan, Zheng, Liang, Yang, Yi, Tian, Qi, Wang, Shengjin, 2018. Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline). In: *The European Conference on Computer Vision (ECCV)*.
- Szegedy, Christian, Liu, Wei, Jia, Yangqing, Sermanet, Pierre, Reed, Scott, Anguelov, Dragomir, Erhan, Dumitru, Vanhoucke, Vincent, Rabinovich, Andrew, 2015. Going deeper with convolutions. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Szegedy, Christian, Vanhoucke, Vincent, Ioffe, Sergey, Shlens, Jon, Wojna, Zbigniew, 2016. Rethinking the inception architecture for computer vision. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Tang, Zheng, Naphade, Milind, Liu, Ming-Yu, Yang, Xiaodong, Birchfield, Stan, Wang, Shuo, Kumar, Ratnesh, Anastasiu, David, Hwang, Jenq-Neng, 2019. CityFlow: A city-scale benchmark for multi-target multi-camera vehicle tracking and re-identification. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Tang, Yi, Wu, Di, Jin, Zhi, Zou, Wenbin, Li, Xia, 2017. Multi-modal metric learning for vehicle re-identification in traffic surveillance environment. In: *2017 IEEE International Conference on Image Processing (ICIP)*. IEEE, pp. 2254–2258.
- Tian, Bin, Tang, Ming, Wang, Fei-Yue, 2015. Vehicle detection grammars with partial occlusion handling for traffic surveillance. *Transp. Res. C* 56, 80–93.
- Tishby, Naftali, Pereira, Fernando C., Bialek, William, 1999. The Information Bottleneck Method. pp. 368–377.

- Varior, Rahul Rama, Shuai, Bing, Lu, Jiwen, Xu, Dong, Wang, Gang, 2016. A siamese long short-term memory architecture for human re-identification. In: Computer Vision – ECCV 2016. Springer International Publishing, pp. 135–153.
- Wang, Xiaogang, 2013. Intelligent multi-camera video surveillance: A review. *Pattern Recognit. Lett.* 34 (1), 3–19, Extracting Semantics from Multi-Spectrum Video.
- Wang, Yan, Wang, Lequn, You, Yurong, Zou, Xu, Chen, Vincent, Li, Serena, Huang, Gao, Hariharan, Bharath, Weinberger, Kilian Q., 2018a. Resource aware person re-identification across multiple resolutions. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- Wang, Cheng, Zhang, Qian, Huang, Chang, Liu, Wenyu, Wang, Xinggang, 2018b. Mancs: A multi-task attentional network with curriculum sampling for person re-identification. In: The European Conference on Computer Vision (ECCV).
- Wei, Longhui, Zhang, Shiliang, Gao, Wen, Tian, Qi, 2018. Person transfer gan to bridge domain gap for person re-identification. In: Computer Vision and Pattern Recognition, IEEE Conference on.
- Wu, A., Zheng, W., Lai, J., 2015. Depth-based person re-identification. In: 2015 3rd IAPR Asian Conference on Pattern Recognition (ACPR). pp. 026–030.
- Wu, A., Zheng, W., Lai, J., 2017a. Robust depth-based person re-identification. *IEEE Trans. Image Process.* 26 (6), 2588–2603.
- Wu, Ancong, Zheng, Wei-Shi, Yu, Hong-Xing, Gong, Shaogang, Lai, Jianhuang, 2017b. Rgb-infrared cross-modality person re-identification. In: The IEEE International Conference on Computer Vision (ICCV).
- Xu, Jing, Zhao, Rui, Zhu, Feng, Wang, Huaming, Ouyang, Wanli, 2018. Attention-aware compositional network for person re-identification. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- Yao, Hantao, Zhang, Shiliang, Zhang, Yongdong, Li, Jintao, Tian, Qi, 2017. Deep representation learning with part loss for person re-identification. *CoRR*, abs/1707.00798.
- Ye, Mang, Liang, Chao, Wang, Zheng, Leng, Qingming, Chen, Jun, 2015. Ranking optimization for person re-identification via similarity and dissimilarity. In: Proceedings of the 23rd ACM International Conference on Multimedia. In: MM '15, ACM, New York, NY, USA, pp. 1239–1242.
- Ye, M., Liang, C., Yu, Y., Wang, Z., Leng, Q., Xiao, C., Chen, J., Hu, R., 2016. Person reidentification via ranking aggregation of similarity pulling and dissimilarity pushing. *IEEE Trans. Multimed.* 18 (12), 2553–2566.
- Yi, D., Lei, Z., Liao, S., Li, S.Z., 2014. Deep metric learning for person re-identification. In: 2014 22nd International Conference on Pattern Recognition. pp. 34–39.
- Yu, Rui, Dou, Zhiyong, Bai, Song, Zhang, Zhaoxiang, Xu, Yongchao, Bai, Xiang, 2018. Hard-aware point-to-set deep metric for person re-identification. In: The European Conference on Computer Vision (ECCV).
- Zhang, Yiheng, Liu, Dong, Zha, Zheng-Jun, 2017. Improving triplet-wise training of convolutional neural network for vehicle re-identification. In: 2017 IEEE International Conference on Multimedia and Expo (ICME). IEEE, pp. 1386–1391.
- Zhao, Fang, Ghorpade, Ajinkya, Pereira, Francisco Câmara, Zegras, Christopher, Ben-Akiva, Moshe, 2015. Stop detection in smartphone-based travel surveys. *Transp. Res. Procedia* 11, 218–226.
- Zhao, Liming, Li, Xi, Zhuang, Yuetong, Wang, Jingdong, 2017. Deeply-learned part-aligned representations for person re-identification. In: The IEEE International Conference on Computer Vision (ICCV).
- Zhao, Junxuan, Xu, Hao, Liu, Hongchao, Wu, Jianqing, Zheng, Yichen, Wu, Dayong, 2019. Detection and tracking of pedestrians and vehicles using roadside lidar sensors. *Transp. Res. C* 100, 68–87.
- Zheng, Liang, Shen, Liyue, Tian, Lu, Wang, Shengjin, Wang, Jingdong, Tian, Qi, 2015. Scalable person re-identification: A benchmark. In: Computer Vision, IEEE International Conference on.
- Zhong, Zhun, Zheng, Liang, Cao, Donglin, Li, Shaozi, 2017a. Re-ranking person re-identification with k-reciprocal encoding. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- Zhong, Zhun, Zheng, Liang, Kang, Guoliang, Li, Shaozi, Yang, Yi, 2017b. Random erasing data augmentation. *CoRR*, abs/1708.04896.
- Zhong, Zhun, Zheng, Liang, Zheng, Zhedong, Li, Shaozi, Yang, Yi, 2018. Camera style adaptation for person re-identification. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- Zhou, Yi, Shao, Ling, 2017. Cross-view gan based vehicle generation for re-identification. In: Tae-Kyun Kim, Gabriel Brostow, Mikolajczyk, Krystian (Eds.), Proceedings of the British Machine Vision Conference (BMVC). BMVA Press, pp. 186.1–186.12.