

UNCERTAINTY AWARE MULTITASK PYRAMID VISION TRANSFORMER FOR UAV-BASED OBJECT RE-IDENTIFICATION

Syeda Nyma Ferdous and Xin Li

West Virginia University
Lane Dept. of Comp. Sci. and Elec. Engr.
Morgantown, WV 26506-6109

Siwei Lyu

State University of New York at Buffalo
Department of Computer Science and Engineering
317 Davis Hall, Buffalo, NY14260

ABSTRACT

Object Re-Identification (ReID), one of the most significant problems in biometrics and surveillance systems, has been extensively studied by image processing and computer vision communities in the past decades. Learning a robust and discriminative feature representation is a crucial challenge for object ReID. The problem is even more challenging in ReID based on Unmanned Aerial Vehicle (UAV) as the images are characterized by continuously varying camera parameters (e.g., view angle, altitude, etc.) of a flying drone. To address this challenge, multiscale feature representation has been considered to characterize images captured from UAV flying at different altitudes. In this work, we propose a multitask learning approach, which employs a new multiscale architecture without convolution, Pyramid Vision Transformer (PVT), as the backbone for UAV-based object ReID. By uncertainty modeling of intraclass variations, our proposed model can be jointly optimized using both uncertainty-aware object ID and camera ID information. Experimental results are reported on PRAI and VRAI, two ReID data sets from aerial surveillance, to verify the effectiveness of our proposed approach.

Index Terms— UAV-based object ReID, Pyramid Vision Transformer, Uncertainty Modeling, Multitask Learning

1. INTRODUCTION

Object Re-Identification (ReID) [1, 2], the task of matching a particular object across different camera views, has been widely studied due to its applications in visual surveillance, especially in the field of person and vehicle tracking. Most of the existing work on object ReID is mainly focused on tackling this problem in a normal surveillance domain, e.g. security cameras installed on the top of a building. With the rapid development in the UAV industry, visual surveillance using UAV devices has received increasing attention, such as normal surveillance. However, ReIDs based on drones or UAVs [3, 4] have remained an under-researched topic. Unlike the normal domain, the ReID of UAV objects is arguably more challenging because drone-based images often contain more



Fig. 1: Example images showing the need for multiscale feature fusion. The images are taken by drones flying at multiple altitudes. A single scale can only capture specific parts of objects, lacking the capability to learn fine scale features. Multiscale features can help a model learn highly discriminative feature space by fusing information across multiple scales.

uncertainty (e.g., view angles, camera distance, and weather conditions) than standard surveillance images.

In the context of the ReID object based on UAVs, there exist several obstacles arising from the increased uncertainty. As shown in Figure 1, both the scale and pose variations are important factors to consider because the altitude changes of the flying drones can be substantial. As the images captured are by UAV flying at different altitudes, a ReID model operating on a single scale or pose cannot create a descriptive feature space that fully characterizes the entire aerial domain. Therefore, it is natural to develop a ReID model capable of learning multiscale and pose-invariant features. Moreover, to detect subtle changes of objects of interest (e.g., a person wearing different colored dresses, different glasses, shoes, etc.), we target a discriminative feature space that is generalizable to both coarse- and fine-grained features.

Today, the convolutional neural network (CNN) is arguably the most popular backbone in terms of designing ReID classifiers. However, the latest advances in deep learning have seen the great success of Transformers [5, 6] in natural language processing and computer vision. This class of convolution-free models can be a better fit for the problem of object ReID using self-attention. Most existing research on object (Persons, Vehicles) ReID has been conducted for normal surveillance. Limited attention has been paid to object ReID in aerial surveillance. Most approaches (e.g., part-based convolutional baseline [7], generative adaptive

alignment network [8], multiple granularity network [9]) follow CNN-based architectures. Transformer [10] is a new trend in the solution of object ReID, establishing it as a strong baseline that beats current state-of-the-art methods.

Motivated by the recent success of the transformer in normal surveillance object ReID [1], we adopt an all-attention transformer-based approach for UAV-based object Re-ID. To address the challenges arising from aerial images, we propose an uncertainty-aware approach that exploits hierarchical feature maps and global channel attention gate for object Re-identification in the UAV domain. Our proposed approach uses a modified version of Pyramid Vision Transformer (PVT) [11] as the backbone, which is a convolution-free architecture. Then we apply spatial attention [12] to these multi-resolution feature maps to put more focus on important features, filtering out irrelevant ones. To incorporate camera information, we add additional head to the PVT model and optimize the network using both camera ID, object ID, and center loss in a joint fashion. Moreover, we normalize the style variance present in different cameras, incorporating Batch Instance Normalization (BIN) [13] in our model. Finally, with the help of the channel attention aggregation gate [14], the model selectively learns the feature maps with higher weights. Inspired by recent work on modeling feature uncertainty for personal Re-ID [15], we propose to model aerial uncertainty by predicting the variance of data as the model output.

PVT-based object recognition with uncertainty estimation is particularly suitable for UAV-based Re-ID of people and vehicles in aerial surveillance and long-range biometrics [16]. Our **technical contributions** can be summarized as follows: **(I)** We seek to explore a modified Pyramid Vision Transformer (PVT) tailored for object re-identification in UAV-based scenarios. The proposed model utilizes multiscale features to re-identify objects for aerial surveillance. **(II)** Spatial attention module helps focus on the relevant information of the feature maps by filtering out noise. The channel attention gate follows an adaptive fusion scheme, which dynamically selects the appropriate feature maps to exploit channel-wise dependencies. **(III)** We train our model considering uncertainty in terms of object identity and camera identity for multitask learning. The model is regularized using Batch Instance Normalization (BIN) to mitigate style variations in multiple cameras. **(IV)** Our proposed framework achieves state-of-the-art performance on two aerial surveillance datasets, PRAI-1581 [3] and VRAI [4], respectively.

2. METHODOLOGY

In this section, we present our multiscale approach for object Re-Identification. An overview of our proposed method is outlined in Figure 2. We propose a multi-task Pyramid Vision Transformer (PVT) [11], a convolution-free backbone designed to learn multiscale feature maps. We use two heads

for object and camera ID recognition, respectively. Batch Instance Normalization (BIN) is incorporated in our model to achieve camera style invariance. To make feature maps spatially aware of the location of important objects, we apply spatial attention [12] to feature maps of different resolutions. Finally, we combine multiscale feature maps in an adaptive way using a Channel Attention Gate [12, 14]. To make the identifier robust to occlusion, we estimate the aleatory uncertainty [17] present in the data while computing the loss for the model.

2.1. Multi-task Pyramid Vision Transformer

Pyramid Vision Transformer (PVT) [11] generates hierarchical feature maps of multiple resolutions. This architecture has four blocks. Each block is responsible for generating a feature map of certain resolution. Each stage consists of a patch embedding layer and a transformer encoder layer. Each transformer encoder layer is composed of a modified attention layer named spatial reduction attention (SRA) and a feedforward layer. SRA is designed to reduce memory cost so that high-resolution feature maps can be processed. First, the input image with a size of $H \times W \times 3$ is separated into $\frac{HW}{4^2}$ patches and the patch size is 4×4 . Then, each patch is flattened through a linear projection and passed through a number of transformer encoders. After that, we get the output feature map of size $\frac{H}{4} \times \frac{W}{4} \times C_1$. Similarly, we obtain the output feature maps of sizes $\frac{H}{8} \times \frac{W}{8} \times C_2$, $\frac{H}{16} \times \frac{W}{16} \times C_3$, and $\frac{H}{32} \times \frac{W}{32} \times C_4$, respectively.

The original PVT uses non-overlapping patches. Inspired by the recent work of TransReID [1], we have generated overlapping patches using a sliding window or shifting operator. The head in PVT is used to identify objects. Differently from the original PVT, we use an additional head for camera ID recognition. By doing this, we try to take advantage of the camera ID labels to fuse camera information into the original model. As the generated feature maps are produced by multiple transformer encoders, the generated feature space is complex, and it can be confusing for the model to extract meaningful features for the ReID task. To obtain a more refined feature map so that the model can learn more informative spatial features, we apply spatial attention to the feature maps of different resolutions. In spatial attention, both MaxPooling and AveragePooling are performed on the channel dimension, and the pooled feature space is concatenated to generate the 2D spatial attention map.

Additionally, we have added Batch Instance Normalization (BIN) [13] to reduce style variations on multiple cameras. The purpose of using BIN is to normalize the style-preserving discriminative features of the ReID objectReID. To fully exploit the functionality of multiscale features, we combine feature maps of multiple resolutions using a global channel attention gate. To tune the channel weights of multiscale feature maps in a dynamic fashion, we use a shared channel at-

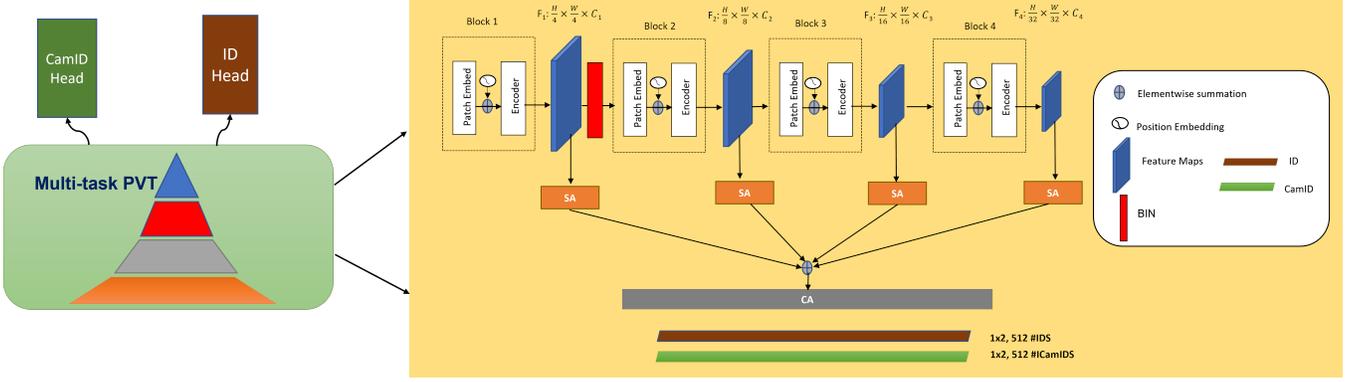


Fig. 2: The overall architecture of our proposed model. The model uses multi-task PVT as the backbone. We apply Spatial Attention (SA) on feature maps of individual scales which helps to focus the network on the most informative features. The channel attention (CA) is a mini-network with shared wights. The Identification (ID) head is responsible for classifying the object instances whereas CamID head gives prediction on camera ID labels considering uncertainty of prediction in both cases. Batch Instance Normalization (BIN) layer helps the network to achieve style generalization among different cameras.

tention gate to learn the inter-channel relationship within a feature map regarding the importance of feature maps. We follow the design procedure mentioned in [14] to implement this gate as a mini-network composed of several global average pooling layers (GAP) and a multi-layer perceptron (MLP) with reduced hidden dimension and one ReLU-activated hidden layer followed by sigmoid activation.

2.2. Loss Function Formulation

To train our model, we use a combination of three uncertainty-aware losses: identity loss, cameraID and center loss. Identity loss consists of softmax cross-entropy loss and triplet loss [18, 19]. Camera ID information is learned through centroid triplet loss [20] and center loss considers the distance from center information.

Uncertainty-aware ID loss: This loss is computed using a hybrid of softmax cross-entropy loss and triplet loss [21] taking into account the uncertainty between classes. The uncertainty-aware softmax cross-entropy loss is calculated by [18]:

$$\mathcal{L}_{softmax} = \frac{1}{2N\sigma(x_i)^2} \sum_{i=1}^N \log(p_{id}(h_i^{id}, y_i)) + \frac{1}{2} \log \sigma(x_i)^2 \quad (1)$$

where N is the number of samples, y_i is the ground truth label, and $\sigma(x_i)^2$ is the variance of the data. We model our uncertainty-aware soft-margin triplet loss using the following formula [19]:

$$\mathcal{L}_{triplet} = \frac{1}{\sigma(x_a)^2} \log(1 + \exp(f(x_a, x_n) - f(x_a, x_p))) + \frac{1}{2} \log \sigma(x_a)^2 \quad (2)$$

where a triplet consists of $\langle x_a, x_p, x_n \rangle$ in which x_a is the anchor image of a person, x_p is the positive anchor image belonging to the identity of the same person, and x_n is the negative anchor image belonging to a different person. Note that triplet loss cannot measure the overall spatial distribution of features, while cross-entropy loss does not have enough discriminant power among features. Therefore, it is better to combine these two as follows:

$$\mathcal{L}_{ua_id} = \mathcal{L}_{softmax} + \mathcal{L}_{triplet} \quad (3)$$

Uncertainty aware camera ID Loss: To tackle intraclass variations arising from view angle, camera style, distance, etc., we apply a soft-margin version of centroid triplet loss since the class centroid can be considered as the mean representation for the retrieval task. Inspired by the unreasonable effectiveness of centroids in image retrieval [20], we propose to calculate uncertainty-aware camera ID loss based on centroids as follows:

$$\mathcal{L}_{ua_camid} = \frac{1}{\sigma(x_a)^2} \log(1 + \exp(f(x_a, c_n) - f(x_a, c_p))) + \frac{1}{2} \log \sigma(x_a)^2, \quad (4)$$

where c_p and c_n are the corresponding centroids of the class for the positive and negative classes.

Uncertainty-aware Center Loss: We also analyze uncertainty-aware center loss [22] using the following formula:

$$\mathcal{L}_{ua_center} = \frac{1}{2\sigma} \sum_{i=1}^B \|f_{t_i} - c_{y_i}\|_2^2, \quad (5)$$

where y_i is the label of the i^{th} image in a mini-batch and B is the batch size. c_{y_i} is the center of deep features in the y_i th class.

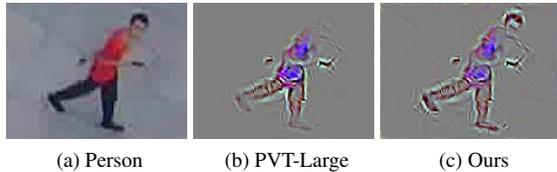


Fig. 3: Visualization of feature maps using guided back-propagation. Baseline PVT-Large [11] fails to retrieve fine features, while ours captures more discriminative features.

The overall loss can be formulated as follows.

$$\mathcal{L}_{total} = \alpha_1 \mathcal{L}_{ua.id} + \alpha_2 \mathcal{L}_{ua.camid} + \alpha_3 \mathcal{L}_{ua.center} \quad (6)$$

Here, α_1 , α_2 , and α_3 are the regularization parameters for the corresponding losses.

3. EXPERIMENTS

3.1. Datasets

We have conducted our experiments on two aerial surveillance datasets named Person ReID for Aerial Imagery (PRAI-1581) [3] dataset and Vehicle Re-identification for Aerial Image (VRAI) [4] dataset. **PRAI** is a newly released aerial surveillance dataset which contains 39,461 person images of 1581 classes captured by two UAV drones with a flight altitude ranging from 20 to 60 meters above the ground. **VRAI dataset** consists around 137,613 images of 13,022 vehicles taken by two UAV drones. This is the largest UAV based vehicle dataset to date.

3.2. Implementation Details

For the PRAI dataset, the training set includes 19,523 images from 782 classes. For the test set, the number of query and gallery images are 4680 and 15258, respectively. For the VRAI dataset, the training set contains 66,113 images with 6,302 classes. For the test set, the query set contains 15,747 images and the gallery set contains 55,753 images, respectively. In our experiment, we investigated PVT, a multiscale transformer, as the backbone network. The backbone network is pre-trained on ImageNet 2012 dataset. We train our model using 4 Titan 1080GTX GPUs. Before training, the images are resized to 224×224 . The batch size is set to 128. ADAM optimizer is used with a momentum of 0.9 and a weight decay of $1e^{-4}$. The learning rate is initialized as 0.000015 with a cosine rate decay. For performance evaluation, we use two metrics: Cumulative Matching Characteristic (CMC) and mean Average Precision (mAP).

3.2.1. Comparison with state-of-the-art

We compare our proposed approach with the latest methods, and the results are reported in Table 1. Our proposed approach outperforms the previous state-of-the-art in both the PRAI-1581 and VRAI dataset, respectively. For the VRAI dataset,

Table 1: Performance comparison with state-of-the-art methods for PRAI-1581 and VRAI.

Method (Person ReID)	Rank-1	mAP
PCB [7, 3]	47.47	37.15
SVDNet [23, 3]	46.10	36.70
MGN [9, 3]	49.64	40.86
OSNet [14, 3]	54.40	42.10
TransReID [1]	56.30	49.81
Ours	59.18	51.45
Method (Vehicle ReID)	Rank-1	mAP
MGN [9, 4]	67.84	69.49
RAM (ResNet-50) [24, 4]	68.58	69.37
RAM (VGG-16) [24, 4]	72.05	57.33
Multi-task+DP [4]	80.30	78.63
TransReID [1]	82.68	81.48
Ours	84.47	82.86

our approach achieves 84.47% Rank-1 accuracy and an mAP of 82.86%. For the PRAI dataset, the accuracy of Rank-1 is 59.18% and the mAP is 51.45%. We report results based on single-query settings for both datasets. It is worth mentioning that the gain over previous SOTA methods is consistent across object domains (person vs. vehicle) and performance metrics (Rank-1 vs. mAP). The performances of our model for mAP and different rank scores are presented in Figure 4 for the PRAI-1581 and VRAI data set. It can be observed that PRAI is a lot more challenging than VRAI because of people’s relatively smaller size, large pose variations, and deformable motion.

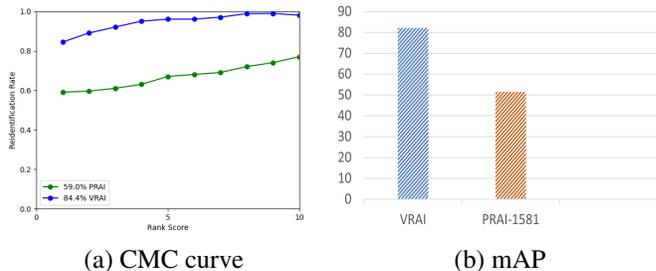


Fig. 4: Performance of our model for PRAI-1581 and VRAI.

4. CONCLUSION

We have presented an uncertainty-aware multiscale transformer-based approach for the UAV-based object Re-ID. Our approach captures the information of instances with different levels of detail by multitasking PVT-based backbone architecture. The proposed model tries to solve object Re-ID as multitask learning problem using a unified framework trained with object ID, camera ID, and center loss. We quantitatively and qualitatively evaluated our proposed method on two UAV-based aerial surveillance datasets. The experimental results demonstrate the superiority of the proposed model over the previous state-of-the-art.

5. REFERENCES

- [1] Shuting He, Hao Luo, Pichao Wang, Fan Wang, Hao Li, and Wei Jiang, “Transreid: Transformer-based object re-identification,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 15013–15022.
- [2] A. Bedagkar-Gala and S. K Shah, “A survey of approaches and trends in person re-identification,” *Image and vision computing*, vol. 32, no. 4, pp. 270–286, 2014.
- [3] Shizhou Zhang, Qi Zhang, Yifei Yang, Xing Wei, Peng Wang, Bingliang Jiao, and Yanning Zhang, “Person re-identification in aerial imagery,” *IEEE Trans. on Multimedia*, vol. 23, pp. 281–291, 2020.
- [4] Peng Wang, Bingliang Jiao, Lu Yang, Yifei Yang, Shizhou Zhang, Wei Wei, and Yanning Zhang, “Vehicle re-identification in aerial imagery: Dataset and approach,” in *ICVV*, 2019, pp. 460–469.
- [5] Ze et al. Liu, “Swin transformer: Hierarchical vision transformer using shifted windows,” *arXiv preprint arXiv:2103.14030*, 2021.
- [6] Daquan et al. Zhou, “Deepvit: Towards deeper vision transformer,” *arXiv preprint arXiv:2103.11886*, 2021.
- [7] Yifan Sun, Liang Zheng, Yi Yang, Qi Tian, and Shengjin Wang, “Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline),” in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 480–496.
- [8] Guan’an Wang, Tianzhu Zhang, Jian Cheng, Si Liu, Yang Yang, and Zengguang Hou, “Rgb-infrared cross-modality person re-identification via joint pixel and feature alignment,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 3623–3632.
- [9] Guanshuo Wang, Yufeng Yuan, Xiong Chen, Jiwei Li, and Xi Zhou, “Learning discriminative features with multiple granularities for person re-identification,” in *Proceedings of the 26th ACM international conference on Multimedia*, 2018, pp. 274–282.
- [10] Alexey et al. Dosovitskiy, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [11] Wenhai et al. Wang, “Pyramid vision transformer: A versatile backbone for dense prediction without convolutions,” *arXiv preprint arXiv:2102.12122*, 2021.
- [12] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon, “Cbam: Convolutional block attention module,” in *ECCV*, 2018, pp. 3–19.
- [13] Hyeonseob Nam and Hyo-Eun Kim, “Batch-instance normalization for adaptively style-invariant neural networks,” *arXiv preprint arXiv:1805.07925*, 2018.
- [14] Kaiyang Zhou, Yongxin Yang, Andrea Cavallaro, and Tao Xiang, “Omni-scale feature learning for person re-identification,” in *ICCV*, 2019, pp. 3702–3712.
- [15] Tianyuan Yu, Da Li, Yongxin Yang, Timothy M Hospedales, and Tao Xiang, “Robust person re-identification by modelling feature uncertainty,” in *ICCV*, 2019, pp. 552–561.
- [16] Mang Ye, Jianbing Shen, Gaojie Lin, Tao Xiang, Ling Shao, and Steven CH Hoi, “Deep learning for person re-identification: A survey and outlook,” *IEEE Trans. on PAMI*, 2021.
- [17] Guotai Wang, Wenqi Li, Michael Aertsen, Jan Deprest, Sébastien Ourselin, and Tom Vercauteren, “Aleatoric uncertainty estimation with test-time augmentation for medical image segmentation with convolutional neural networks,” *Neurocomputing*, vol. 338, pp. 34–45, 2019.
- [18] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, Zhibo Chen, and Shif-Fu Chang, “Rethinking classification loss designs for person re-identification with a unified view,” *arXiv preprint arXiv:2006.04991*, 2020.
- [19] Alexander Hermans, Lucas Beyer, and Bastian Leibe, “In defense of the triplet loss for person re-identification,” *arXiv preprint arXiv:1703.07737*, 2017.
- [20] Mikolaj Wicczorek, Barbara Rychalska, and Jacek Dabrowski, “On the unreasonable effectiveness of centroids in image retrieval,” *arXiv preprint arXiv:2104.13643*, 2021.
- [21] Zihao Hu, Huiyan Wu, Shengcai Liao, Hai-Miao Hu, Si Liu, and Bo Li, “Person re-identification with hybrid loss and hard triplets mining,” in *2018 IEEE Int. Conf. on Multimedia Big Data (BigMM)*, 2018, pp. 1–5.
- [22] Hao Luo, Youzhi Gu, Xingyu Liao, Shenqi Lai, and Wei Jiang, “Bag of tricks and a strong baseline for deep person re-identification,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 2019, pp. 0–0.
- [23] Yifan Sun, Liang Zheng, Weijian Deng, and Shengjin Wang, “Svdnet for pedestrian retrieval,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3800–3808.
- [24] Xiaobin Liu, Shiliang Zhang, Qingming Huang, and Wen Gao, “Ram: a region-aware deep model for vehicle re-identification,” in *ICME*, 2018, pp. 1–6.