

# SUES-200: A Multi-height Multi-scene Cross-view Image Benchmark Across Drone and Satellite

Runzhe Zhu, Ling Yin, Mingze Yang, Fei Wu, Yuncheng Yang, Wenbo Hu

**Abstract**—Cross-view image matching aims to match images of the same target scene acquired from different platforms. With the rapid development of drone technology, cross-view matching by neural network models has been a widely accepted choice for drone position or navigation. However, existing public datasets do not include images obtained by drones at different heights, and the types of scenes are relatively homogeneous, which yields issues in assessing a model’s capability to adapt to complex and changing scenes. In this end, we present a new cross-view dataset called SUES-200 to address these issues. SUES-200 contains 24120 images acquired by the drone at four different heights and corresponding satellite view images of the same target scene. To the best of our knowledge, SUES-200 is the first public dataset that considers the differences generated in aerial photography captured by drones flying at different heights. In addition, we developed an evaluation for efficient training, testing and evaluation of cross-view matching models, under which we comprehensively analyze the performance of nine architectures. Then, we propose a robust baseline model for use with SUES-200. Experimental results show that SUES-200 can help the model to learn highly discriminative features of the height of the drone.

**Index Terms**—Cross-view Image Matching, Drone, Benchmark, Image Retrieval, Pipeline, Geo-localization

## I. INTRODUCTION

CROSS-VIEW matching [1] is an essential topic in computer vision research. This technique can be applied in many domains, such as localization, navigation, autonomous driving, and object detection. Satellite and drone platforms are the primary sources of images used as input for the matching. A standard cross-view matching system works as follows. Given an image to be retrieved in a query dataset from one view, the matching system finds an image under the exact location in a large-scale candidate (gallery) dataset of another view. There are two main tasks:

Task 1: Drone view target localization (Drone → Satellite) and Task 2: Drone navigation (Satellite → Drone).

Thus, the key to the effectiveness of the matching techniques is learning the discriminative features of images that are invariant under different views.

Previous studies on cross-view matching [2]–[7] mainly focus on the matching between street view and satellite view, or between street view and bird-view. For example, the datasets

Runzhe Zhu, Ling Yin, Mingze Yang, Fei Wu, Yuncheng Yang are with the School of Electronic and Electrical Engineer, Shanghai University of Engineering Science, Shanghai 201602, China (email: m025120503@sues.edu.cn; lyin@sues.edu.cn; ymz871500142@163.com; fei\_wu1@163.com; shawn.yangyc@foxmail.com. ) Corresponding author: Fei Wu.

Wenbo Hu is with the School of Communication and Information Engineering, Shanghai University, Shanghai 200444, China (email: wenbohu@shu.edu.cn. )

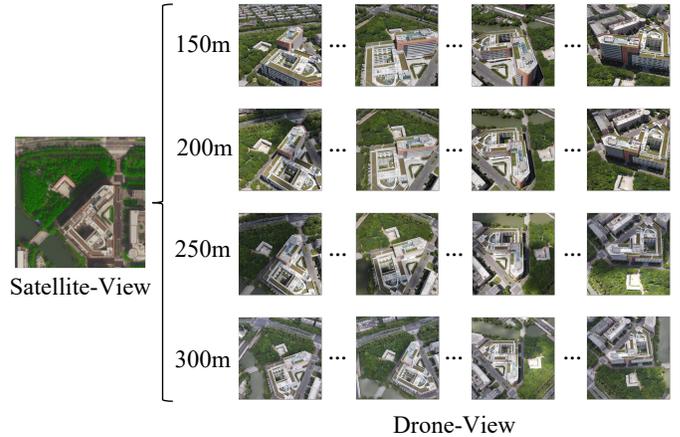


Fig. 1. A representative target scene of SUES-200 contains fifty drone view images from four heights and one satellite view image.

CVUSA [8] and CVACT [9] used panoramic street view and satellite views of the same target scene to construct a cross-view image pair for training a deep neural network model. The quality of matching between street view and satellite images is limited by the much smaller spatial scale of a street view. Thus, street view tends to be obscured and interfered with, resulting in features not properly extracted by models.

With the wide application of drone technology [10]–[12], more and more researchers have been using drone platforms to capture target scenes at different spatial and temporal scales. Traditionally, image matching of drone view and satellite views is relatively limited in the military field; fixed-wing drones are conventionally designed to fly at a specific height and collect images in real-time [13]–[15]. Matching systems are used to match the images captured by a drone with satellite images to infer the drone’s location. This autonomous locating system is not affected by the external environment and has strong robustness in complex electromagnetic environments. Recently, rotary-wing drones have been gained wide applications. How to use such vehicles for positioning in low airspace has become a hot research issue.

Recently, new progress has been made in cross-view view matching research. Zheng et al. [16] established the first drone-based multi-source cross-view matching dataset, namely University-1652, which contains images from three perspectives, including street view, aerial drones, and satellite. They also published a baseline for multi-branch CNN networks. [17]–[21], and matching accuracy was significantly improved

in a more in-depth study. However, this dataset still involves a few problems. For example, only synthetic images of drone views are included, which lack realistic variations in lighting. Similarly, differences in images captured by drones at different heights are not distinguished. Moreover, the captured scenes are of a single type, mostly buildings on campuses. These problems limit the ability of learning models trained on this dataset to differentiate different types of scenarios. Furthermore, such models are unable to extract robust features from images captured at low heights.

To address these problems, we propose a multi-height, multi-scene dataset including images from both drones and satellites based on the University-1652 dataset, called SUES-200. SUES-200 contains a wider variety of scenes, such as parks, schools, lakes, and public buildings. For each scene, data collected at four different heights (150m, 200m, 250m, and 300m). All of the included images were recorded from onboard drones in flight in real world. SUES-200 contains 200 target scenes, 120 of which are specified for use as a training set, and 80 scenes of which are designated as a testing set. Some samples in SUES-200 are shown in Figure 1.

Traditional evaluation metrics for cross-view matching datasets are Recall@K [22] and AP. However, these measures are not suitable for the characteristics of our new SUES-200 dataset, because the differences in drone views at different heights are not taken into account. Moreover, drones encounter diverse interference when flying outdoors. Therefore, we developed a new evaluation system that focuses on three aspects of the model, including 1) robustness at different heights; 2) robustness to uncertainties; and 3) inference speed; In addition, we provide a pipeline dedicated to cross-view matching, which helps to improve the efficiency of training, testing, and model evaluation.

As an experiment, we train and test feature extractors of different deep neural network (DNN) architectures on SUES-200 using the pipeline developed in this work. The model with the best overall evaluation results is released as the baseline model of SUES-200. We also evaluate the effects of multi-angle feature fusion on matching results and compare the performance of transferred learning models. We perform ablation studies to evaluate each component of the baseline model. Our results show that SUES-200 can help neural models learn high-level features in various scenes captures from different heights. With increasing height, drone footage is gradually less affected by the environment and camera pose and achieves better performance metrics.

We release a ViT-based model as the baseline of SUES-200. For Drone  $\rightarrow$  Satellite, baseline achieves 59.32, 62.30, 71.35, and 77.17 Recall@1 accuracy at heights of 150m, 200m, 250m, and 300m, respectively. For Satellite  $\rightarrow$  Drone, baseline achieves 82.50, 85.00, 88.75, and 96.25 Recall@1 accuracy in 150m, 200m, 250m, and 300m, respectively. This baseline also showed strong robustness to different heights and uncertainties. ViT is a very general scheme compared to other CNN-based algorithms, although its computational complexity is large. The entire dataset, as well as the code for the evaluation, is available at <https://github.com/Reza-Zhu/SUES-200-Benchmark>.

The main contributions of this study are summarized as follows.

- We build a new cross-view matching dataset: SUES-200, which provides diverse scenes and height views for each scene. All images are acquired in real environments of multiple types of scenes, including real-world light, shadow transformations and disturbances. The datasets, as well as the code for the evaluation, are available at <https://github.com/Reza-Zhu/SUES-200-Benchmark>.
- We propose a new evaluation system based on the characteristics of SUES-200 to evaluate the robustness of matching models for different heights, robustness to uncertainties, and the speed at which inferences are performed, together with to the classical Recall@K and AP.
- We establish an efficient pipeline to train and test different matching models and release the baseline model of SUES-200 according to the comprehensive evaluation results.

## II. RELATED WORK

### A. Cross-view Datasets

Previous cross-view datasets mostly focused on images collected from the same location with different viewpoints via different platforms such as panoramic cameras, satellites, drones, and smartphones. For example, the dataset [2] comprises publicly available data, containing a total of 78K data pairs. Each pair consists of two views, including an aerial or bird's-eye view, and the other view is the street view. Tian *et al.* [4] collected images from several locations in a city and constructed image pairs with bird's-eye and street views. Tian incorporates semantic information to label buildings in images from different views and contains an object detection module in its network structure. The experimental results were evaluated in terms of PR curves and AP. CVUSA [23] is a standard cross-view dataset, consisting of image pairs of panoramic street views and satellite views. CVACT [9] is a larger panoramic dataset with improved satellite image resolution and more testing sets. Moreover, GPS tags to scenes are supplemented. Both CVACT and CVUSA use Recall@K to evaluate matching results. University-1652 was proposed by Zheng for multi-source cross-view scene matching *et al.* [16] as the first geo-localization dataset based on drone footage. It contains image data triplets with satellite, drone, and street views for 1652 buildings in 72 universities. University-1652 generally includes one satellite image, fifty-four images captured by aerial drones, and multiple street view images for a given location. Due to the expense of real-world flight, the drone data in this database was obtained by simulated flights in Google Earth. The drone simulation flight route circles around the target scene and gradually drops in height. University-1652 uses Recall@K and AP to evaluate matching results. Inspired by University-1652, we constructed the SUES-200 dataset to emphasize differences in images acquired by drones at different heights. In addition, we extended the types of scenes, all of which were captured in real scenarios.

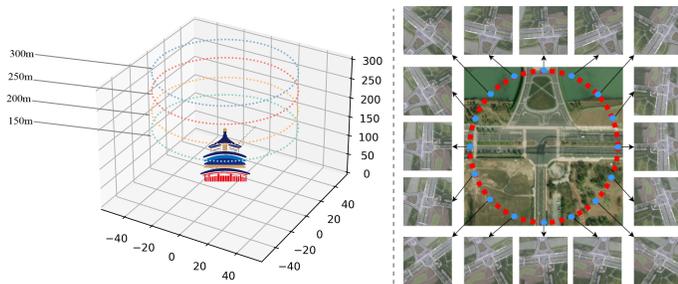


Fig. 2. The flight height of the drone when collecting images is 150m, 200m, 250m, and 300m. The flight trajectory is one circle around the target scene

TABLE I  
OVERVIEW OF CROSS-VIEW APPROACHES.

Approach	Feature Extractor	Author
CVUSA	VGG16	Workman <i>et al.</i> [23]
CVACT	7-layer CNN	Liu <i>et al.</i> [9]
University-1652	ResNet-50	Zheng <i>et al.</i> [16].
LCM	ResNet-50	Ding <i>et al.</i> [17]
LPN	ResNet-50	Wang <i>et al.</i> [19]
PCL	ResNet-50	Tian <i>et al.</i> [20]
MSBA	ResNet-50	Zhuang <i>et al.</i> [18]

### B. Cross-view Methods

Traditional cross-view matching methods [24]–[26] are based on hand-crafted feature descriptors such as SIFT [24], SURF [25], and ORB [27]. However, these feature extraction methods are not robust and are susceptible to uncertainties such as lighting and occlusion, especially for drones flying at heights. False or missing matches typically occur frequently due to excessive differences between the acquired images and the satellite view images. Since the publication of the University-1652 dataset, considerable progress has been made in the past years in deep learning methods. Liu *et al.* [17] proposed LCM, which utilized ResNet [28] as a backbone network and trained the image retrieval problem as a classification problem. The LCM improved the Recall@1 and AP by 5-10 % over the baseline of University-1652. Wang *et al.* [19] designed LPN to consider the contextual information of neighboring regions. The LPN used a square-ring partition strategy to divide feature maps, which provided good robustness to changes in rotation. LPN achieved good performance on University-1652, CVACT, and CVUSA. Tian *et al.* [20] presented a method that integrated the spatial correspondence between the satellite views and information on the surrounding area. This was performed in two steps, first converting the tilted view of the drone into a vertical view by perspective transformation and then transforming the image of the drone view to be closer to the satellite view using a conditional GAN [29]. The experimental results show that this method improved accuracy by 5% over LPN on University-1652. Inspired by the development of attention mechanisms, Zhuang *et al.* [18] developed MSBA to eliminate the differences in images acquired from different viewpoints. MSBA cuts an image into several parts with different scales. Based on that division, a self-attention mechanism is used for more effective feature extraction. They showed that MSBA performed better

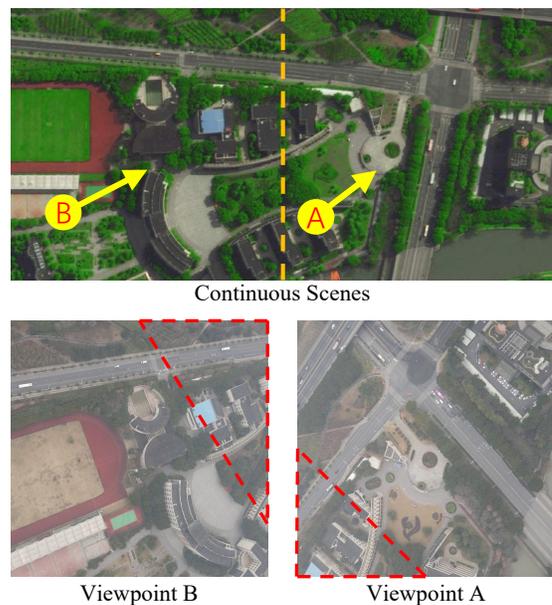


Fig. 3. Continuous Scenes display two continuous satellite view images. Yellow arrows indicate the directions of the drones' viewpoints. Two corresponding drone view images are shown below, with the area overlapping in the images marked by the dashed red line.

than LPN in terms of accuracy and inference efficiency. Table I gives an overview of existing approaches. Importantly, most approaches adopt the same backbone network and were not tested for other feature extractors. In contrast, in the present work, we trained several cross-view matching models and then tested and evaluated the performance of different backbone networks such as VGG [30], ResNet [28], and DenseNet [31] in extracting features at different heights using the pipeline.

## III. SUES-200 DATASET

### A. Dataset Description

SUES-200 is a cross-view matching dataset with the characteristics of multiple sources, multiple scenes, and panoramic views. We collected multi-source images of satellite views and corresponding drone views at 200 locations around the Shanghai University of Engineering and Science(SUES). We used 0001, 0002, ... 0200 to distinguish the images obtained in different scenes, and numbers 1-200 represent specific scenes. Specifically, to enable the model to learn highly discriminative features at different heights, we collected drone view images at 150m, 200m, 250m, and 300m. SUES-200 includes a broader range of scene types, not limited to campus buildings, containing parks, schools, lakes, and public buildings. The rich multi-type scenes enable the models to learn features that can be adapted to real environments.

Satellite-view images are obtained from AutoNavi Map and Bing Maps. A single satellite image is included for each location. The schematic diagram of the drone flight is shown in Figure 2. The drone flight path was set to a curve in different heights to capture multi-angle information of target scenes. We sampled 50 frames uniformly from the flight video recorded by

TABLE II  
COMPARISON BETWEEN SUES-200 AND OTHER CROSS-VIEW DATASETS.

Datasets	SUES-200	University-1652 [16]	CVUSA [23]	Tian et al [4].
Platform	Drone, Satellite	Drone, Ground, Satellite	Ground, Satellite	Ground, 45° Aerial
Target	Diversity	Building	User	User
Height difference	TRUE	FALSE	FALSE	FALSE
Training	120 * 51	701 * 71.64	35.5k * 2	15.7k * 2
Images/Location	50 + 1	51 + 16.64 + 1	1 + 1	1 + 1
Evaluation	Recall@K & AP & Robustness & Inference Speed	Recall@K & AP	Recall@K	PR&AP

TABLE III  
STATISTICS OF SUES-200 TRAINING AND TEST SETS, INCLUDING THE IMAGE NUMBER AND THE SCENE NUMBER OF TRAINING SET, TESTING SET.

Training Dataset			
Views	Locations	Images at Each Height	Total
Drone	120	6000	24000
Satellite	120	–	120
Testing Dataset			
Views	Locations	Images at Each Height	Total
Drone query	80	4000	16000
Satellite query	80	–	80
Drone gallery	200	10000	40000
Satellite gallery	200	–	200

the drone. Overall, every location includes one satellite view image and 50 drone view images.

In addition, multiple satellite images are consecutively selected in the same area by SUES-200. When the drone flies in one of the locations, the image includes information about nearby locations. As shown in Figure 3, there is some overlap between drone maps of different scenarios. It is desirable that the cross-view matching models could pay attention to the main feature in the scene without the effect of overlapping regions.

In order to prevent information loss due to image resolution, drone images in SUES-200 use the original resolution of  $1080 \times 1080$  and satellite images use the resolution of  $512 \times 512$ . The dataset includes 200 locations with 50 drone images and 1 corresponding satellite image for each location. SUES-200 is divided into training and testing sets, with 120 locations designated for training and 80 locations as testing data. To accomplish the two tasks mentioned in the introduction, the testing data include the query drone dataset, query satellite dataset, gallery drone dataset, and gallery satellite dataset. Among these, the gallery dataset contains the testing data and adds the training data as confusion data to increase the difficulty of matching.

In the testing phase, we consider Task 1 and Task 2 as image retrieval tasks. Taking Drone  $\rightarrow$  Satellite as an example, the query is a drone image, the gallery is satellite image. The model first extracts features from the images in the gallery set and stores them locally. Then, a single query image is fed into the model to extract features and calculate the distance between the query feature and gallery images. The image pair with the shortest distance between the drone view image and the satellite view image is considered the matching result. Some statistics on the datasets are shown in Table II and Table III.

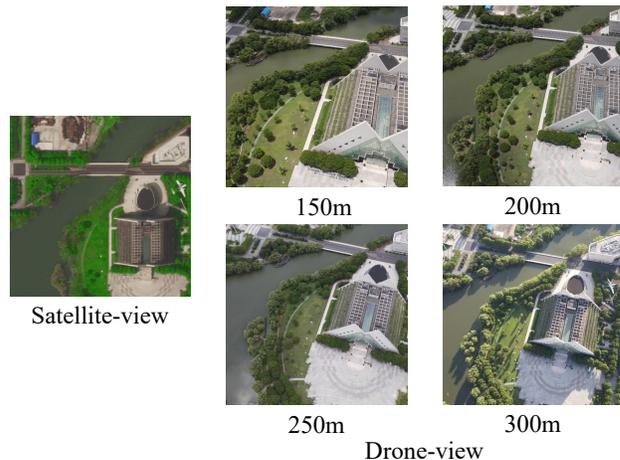


Fig. 4. As the height rises, the images captured from the perspective of the drone become increasingly similar to the satellite view.

Finally, we summarize the new characteristics of the SUES-200 dataset.

- 1) **Multi-height:** SUES-200 contains data collected at different heights: 150m, 200m, 250m, and 300m, and can evaluate model metrics at different heights. To the best of our knowledge, SUES-200 is the first cross-view dataset to include images recorded by cameras on drone vehicles flying at different heights.
- 2) **Multi-scene:** SUES-200 contains data from different types of scenes. This can help models extract invariant features in more scenes and expand the scope of scene applications for drone-based cross-view matching techniques.
- 3) **Continuous-scenes:** Some of the included target scenes were collected in the same area. The drone image is thus affected by the surrounding scene; for example, information from another scene may be recorded. This poses a challenge to the ability of trained models to differentiate between scenes, but this is also realistic for practical application environments.

## B. Evaluation Protocol

In this subsection, we introduce the evaluation system of SUES-200. In response to the existing real-world problems, in addition to the traditional Recall@K [3], [9], [32] and AP [4], [33] evaluation metrics, we propose a method to measure model robustness at different heights, as well as a method to

measure the robustness of trained models to uncertainties and a method to evaluate inference speed.

**Recall@K and AP.** SUES-200 contains 200 target scenes, including 120 scenes for training and 80 scenes for testing. Among these, 120 scenes from the training set are also included in the gallery as distractors. There is no overlap between the training and testing data. Recall@K (R@K) represents the probability that a correct match appears in the top-k ranked retrieval results. Recall@1 is very sensitive to the position of the first true-matched image appearing in the ranking of the matching result. A higher recall score shows a better performance of the network. The AP is the area under the precision-recall(PR) curve, which considers the position of all true-matched images in the evaluation. Recall@K is defined as follows.

$$\text{Recall@K} = \begin{cases} 1, & \text{if } \text{order}_{\text{true}} < K + 1 \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

AP is formulated as follows:

$$\text{AP} = \frac{1}{m} \sum_{h=1}^m \frac{p_{h-1} + p_h}{2}, \text{ where } p_0 = 1 \quad (2)$$

$$p_h = \frac{T_h + 1}{T_h + F_h} \quad (3)$$

where  $m$  is the number of true-matched images for a query,  $T_h$  and  $F_h$  are the numbers of true-matched images and false-matched images before the  $(i + 1)$ -th true-matched image in the matching.

**Robustness at different heights.** SUES-200 differentiates the images acquired by drones at different heights, as shown in Figure 4. Measuring the robustness of the model at different heights is also an important evaluation index. The drone images appear most similar to satellite view images at 300m. As the height decreases, the drone's field of view gradually narrows. The size of the target scene becomes larger, and more detailed information is presented. These factors make it increasingly difficult for the model to distinguish the variations across scenarios. To evaluate the model's robustness to height, we set the Recall@1 at 300m as the baseline. Then, the Recall@1 at other heights is divided by the baseline to evaluate the reduction in accuracy with height, which is calculated as follows.

$$\text{RDR}_n = 1 - \frac{R@1_n}{R@1_{300m}} \quad (4)$$

$\text{RDR}_n$  represents the recall degradation rate (RDR) of the model at a given height  $n$ . The RDR of the model between 300m to 150m directly reflects the robustness of the height of the model. Finally, we denote the overall robustness protocol of the height as:

$$\text{RDR}_{150m} = 1 - \frac{R@1_{150m}}{R@1_{300m}} \quad (5)$$

The larger the  $\text{RDR}_{150m}$ , the less robust the model is to height, and vice versa.

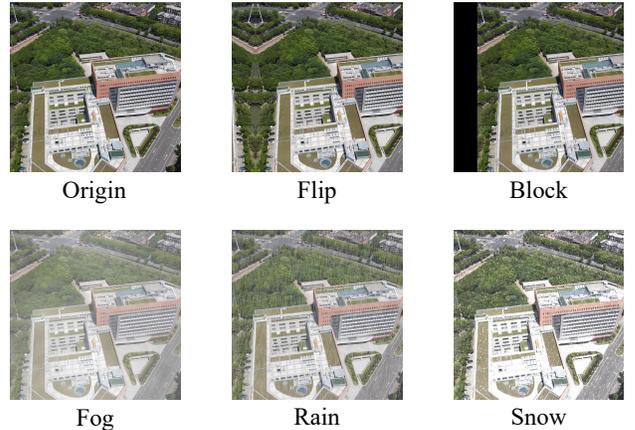


Fig. 5. We list five uncertainties including flip, block, fog, rain, and snow.

**Robustness to uncertainties.** In practice, the images captured by drone vehicles are often disturbed by various uncertainties, such as the target being obscured or offset and various weather factors. To evaluate the performance of the model under those uncertainties, we simulated these effects by applying augmentation to the drone images in the query set. We considered five types of factors, including flip, block, fog, snow, and rain, in Figure 5. We denote the original AP of the model as  $AP_{\text{origin}}$ , and the disturbed AP at a certain height as  $AP_i$ . The  $n = 4$  indicates four times this height. We calculate the average rate of degradation of precision (RDP) for a given model at four heights to indicate its robustness. The equation is shown as follows.

$$\text{RDP} = \frac{\sum_{i=1}^{n=4} 1 - AP_i / AP_{\text{origin}}}{4} \quad (6)$$

**Inference Speed.** In the actual application process, the inference speed of the model is a significant concern. Therefore, we refer to the formula mode of [18] to evaluate the inference speed. [18] proposed a "real-time" method to evaluate the inference time of a single query. However, we considered that this approach is not sufficiently comprehensive to evaluate the overall performance because its definition of "real-time" only focuses on the inference time of the query image. Therefore, we present inference speed to assess the combined inference time of the query and gallery images. We chose a base model with relatively fast inference performance and take its inference time as the benchmark speed. The benchmark is set as 1.00, and other inference speeds are then denoted as  $1.00 \times T (0 < T < +\infty)$ .

## IV. METHOD

### A. Pipeline

We established a pipeline to solve the cross-view matching problem. This approach provides a standardized method to efficiently train, test, and evaluate different models. As shown in Figure 6, in this pipeline, the input is the cross-view matching dataset, and the output is the values of each

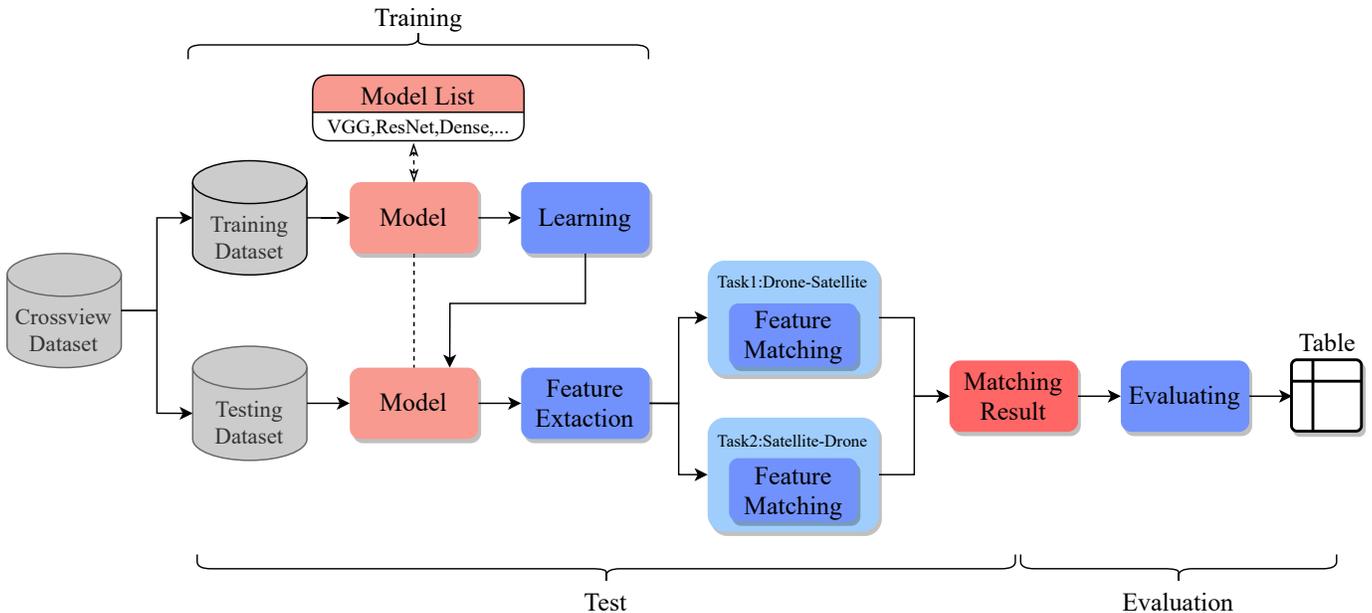


Fig. 6. The pipeline reads images from the dataset and sends them to the currently selected model for training. After training, the model with the best parameters is selected and sent to Task 1 or Task 2 for testing, and the evaluation module evaluates the test results to form an evaluation table.

evaluation index, in which the model is a deep neural network constructed by the user. The network is divided into a backbone network and a classification network part. Selecting different feature extractors in the “Model List” replaces the backbone network part in the corresponding network structure, and the user can also customize its network structure. Details of the network structure of deep neural networks are provided in the next section. The images in the testing set are input to the model, which extracts features and completes Task 1: Drone  $\rightarrow$  Satellite, and Task 2: Satellite  $\rightarrow$  Drone, and the obtained feature matching results are finally passed to the evaluation unit to obtain the evaluation table.

### B. Network Architecture and Loss Function

The drone and satellite images included in SUES-200 originate from different sources, but there are still some similarities. Our deep learning network extracts robust and invariant features separately and maps them to a high-dimensional space for the following matching task. After referring to previous studies, we constructed a two-branch deep neural network (DNN) architecture. One branch extracts feature from satellite view images, and the other branch extracts features from drone view images. To test the performance of different DNN structures on different source image feature extractors, we apply network structures that extract features in two branches of backbone networks that are replaceable. Subsequently, we add a shared weight fully connected (FC) layer to unify the feature dimensions. In the training process, we add an MLP block, including a drop-out layer, an FC layer, and a softmax layer, at the end of the branch to treat the processing as a classification task. Each target location is treated as a class to train the entire network. In the testing process, the feature map is unified by FC1, and then the distance between each

feature is calculated by distance measurement algorithms. The architecture of the network is shown in Figure 7.

In recent years, different DNN structures have been extensively developed. ResNet [28] is widely used as a CNN-based backbone network [17]–[20] for feature extraction in the field of cross-view matching due to its clever design structure and excellent performance. With further research on ResNet and the emergence of attention mechanism, ResNet has been further improved to produce extended variants such as SE-ResNet [34], ResNeSt [35], CMAB-ResNet [36], and these models have achieved excellent performance on image classification datasets such as ImageNet [37]. In addition to ResNet, other structured CNNs are also a topic of considerable research interest in recent years, including DenseNet [31], EfficientNet [38], Inception [39]. Moreover, ViT [40] architecture has achieved great success in various computer vision tasks. Is there a more proper feature extractor than ResNet in cross-view matching? In our experiment, we tested the improved CNN-based and other structures on SUES-200 and evaluated these models according to the evaluation system.

For the loss function, since the model training process is considered a multi-classification task, we adopted a cross-entropy as the loss function. Cross entropy is mainly used to determine how close the actual output is to the expected output, i.e., the smaller the cross entropy between the network output and the labels, the better the classification ability of the network.  $z_j^i(y)$  is the logarithm of ground-truth  $y$ , and  $\hat{p}(y|x_j^i)$  is the probability of the predicted outcome of the model equal to ground-truth  $y$ . The mathematical formula is given as follows.

$$\hat{p}(y|x_j^i) = \frac{\exp(z_j^i(y))}{\sum_{c=1}^C \exp(z_j^i(c))} \quad (7)$$

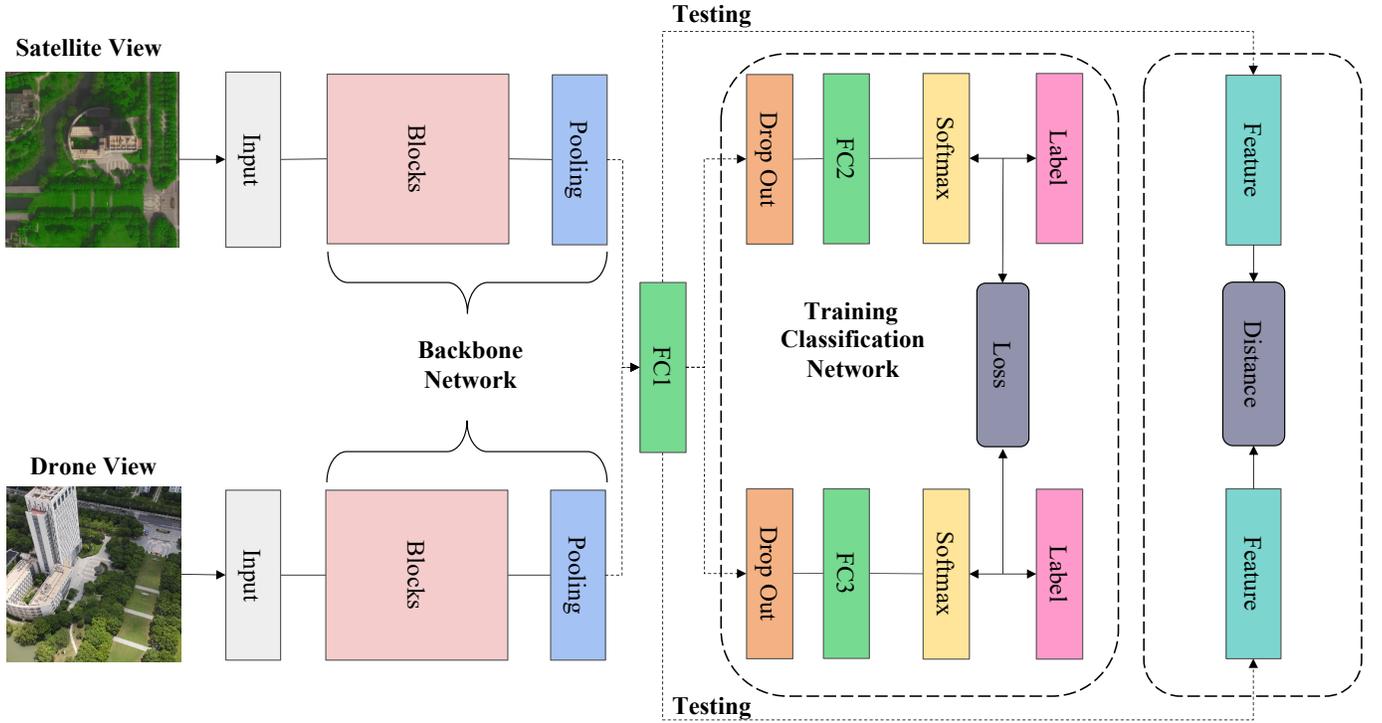


Fig. 7. Basic network architectures for cross-view matching. We apply two-branch network structures with cross-entropy loss to train the model. The feature map is extracted by the backbone network and pooling layer. Then, a shared weight FC layer(FC1) unifies the feature dimension. Next, the feature map is fed into the classifier network for training. In addition, the cosine distance is used to calculate the similarity between the query and candidate images in the gallery for testing.

$$\text{Loss} = \sum_{i,j} -\log(\hat{p}(y|x_j^i)) \quad (8)$$

In the two-branch DNN, both outputs of the model need to be compared with the label to obtain two loss values. Let the loss of drone view be  $L_d$ , and the loss of satellite view be  $L_s$ , and these two loss values are added to get  $L_{total}$ . We optimize the whole network through  $L_{total}$ .

$$L_{total} = L_s + L_d, \quad (9)$$

the query images in the test set are from a drone view and a satellite view. We feed the query images to the model with fixed parameters, remove the classification network from the training layer, and use the backbone network to output the feature vectors directly. The feature vector of the drone-view image is represented as  $f_d$ , and the feature vector of the satellite view is defined as  $f_s$ . Our test aim is to find the most similar set of feature vectors by cosine distance to measure the similarity between  $f_d$  and  $f_s$ .  $f_{di}$  and  $f_{si}$  are parts of the feature vector, and a smaller cosine distance means that the set of features is less similar. A larger cosine distance implies that this pair of feature maps are more similar to each other.

The formula is given as follows.

$$\text{Cosine} = \frac{f_d f_s}{\|f_d\| \times \|f_s\|} = \frac{\sum_{i=1}^n f_{di} f_{si}}{\sqrt{\sum_{i=1}^n (f_{di}^2)} \sqrt{\sum_{i=1}^n (f_{si}^2)}} \quad (10)$$

## V. EXPERIMENT

In this section, we first describe the experimental setup and details, followed by a comprehensive evaluation of multiple feature extractors through the pipeline. The impact of multiple queries on the matching performance is explored. In addition, we test the performance of the transfer learning model on SUES-200. Finally, we implement some classical cross-view matching models on SUES-200.

### A. Implementation Details

Different feature extractors are used in our backbone network, and all of them are loaded with ImageNet's pre-trained weights to speed up the convergence of the model. However, the amount of work required to tune so many models to the optimum is considerable. For training, we applied the grid search method to search for the best learning rate, dropout rate, and weight decay hyperparameter values. The image size is resized to (384, 384) before feeding to the network, and only the basic image augmentation methods are used,

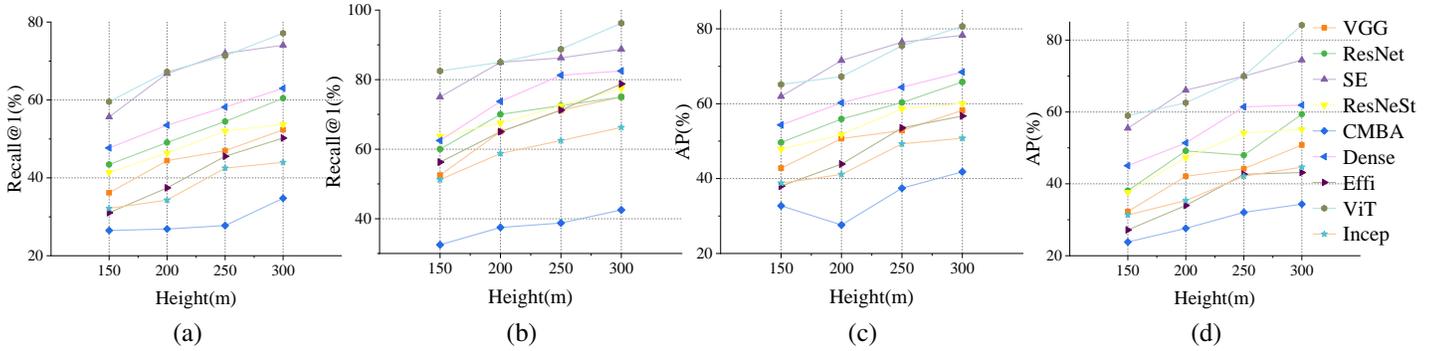


Fig. 8. The Recall@K accuracy curve and AP value curve at 150m, 200m, 250m, and 300m. (a): Recall@1 curve of Drone  $\rightarrow$  Satellite. (b): Recall@1 curve of Satellite  $\rightarrow$  Drone. (c): AP curve of Drone  $\rightarrow$  Satellite. (d): AP curve of Satellite  $\rightarrow$  Drone.

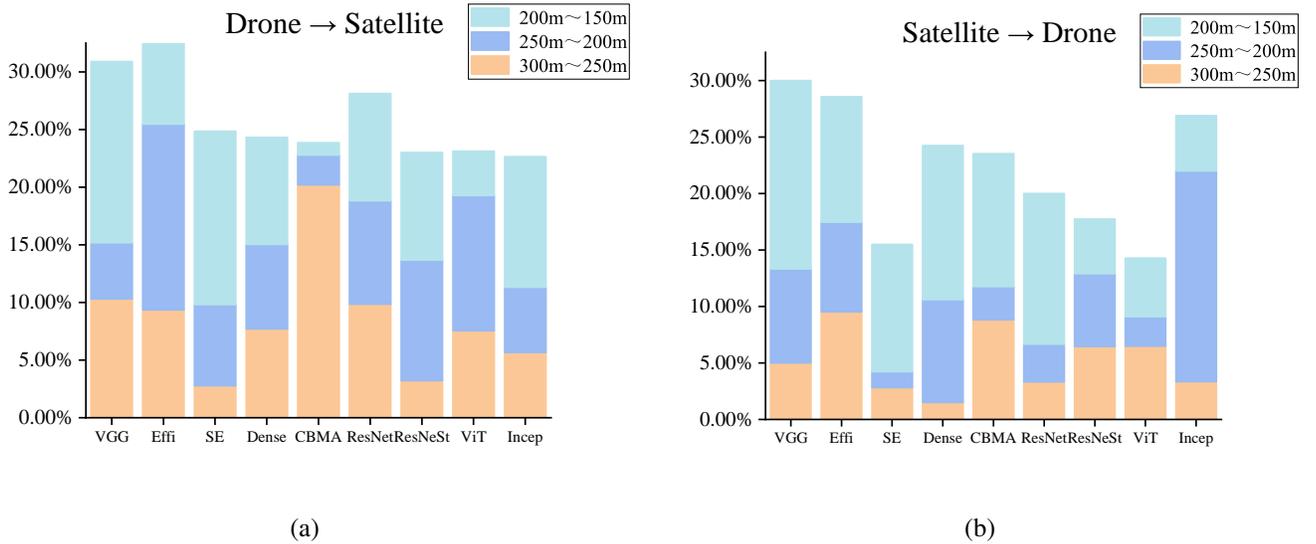


Fig. 9. The robustness of different backbone networks at different Heights. The height of the bars in the bar chart represents the total recall accuracy loss over the height from 300m to 150m. Each of the three colors indicates the loss of accuracy in the respective height interval. (a): Drone  $\rightarrow$  Satellite. (b): Satellite  $\rightarrow$  Drone.

including random cropping and random horizontal flipping. The optimizer of the neural network is SGD (momentum=0.9), and the initial learning rates of the backbone network and the classification network are set to 0.1 times and 1 time of the learning rate. The learning rate decay is MultiStepLR, and the parameters of the classification network are initialized with Kaiming Initialization [41]. In the testing stage, we apply *imgaug* [42] to simulate the unfavorable elements for drone view images. Our model was constructed using the PyTorch framework, and all experiments were conducted on an NVIDIA RTX TITAN GPU.

### B. Evaluation of Different Extractors

We aim to determine whether SUES-200 can help a model learn highly discriminative features, and whether the pipeline could perform the tasks of training, testing, and evaluation efficiently. In this section, we describe experiments conducted to comprehensively evaluate feature extractors of different DNN architectures and use the model with the best experimental results as the baseline model of SUES-200.

**Recall and AP.** Using the pipeline, we quickly train the models on the SUES-200 for testing and evaluation. As shown in Figure 8, we compare the feature extraction capability of different backbone networks by Recall@K and AP. In the drone-view target localization task (Drone  $\rightarrow$  Satellite), ViT achieves 59.57%, 62.30%, 71.35%, 77.15% Recall@1 in four heights, respectively. In the drone navigation task (Satellite  $\rightarrow$  Drone), ViT achieves 82.50%, 85.00%, 88.75%, 96.25% Recall@1 in four heights, respectively. The performance surpasses the results of other common feature extractors such as ResNet. ViT shows considerable potential for cross-view matching as a general framework. The results also show that as the height of the drone increases, the images captured are less affected by the surrounding environment and the camera field view. The images camera acquired by the camera are more similar to satellite images, and the Recall@K and AP of the model are improved.

**Robustness at different heights.** Because the height of the drone affects the accuracy of the matching system, we evaluate the robustness of the model to different heights in positioning or navigation tasks using Equations (4)-(5). As

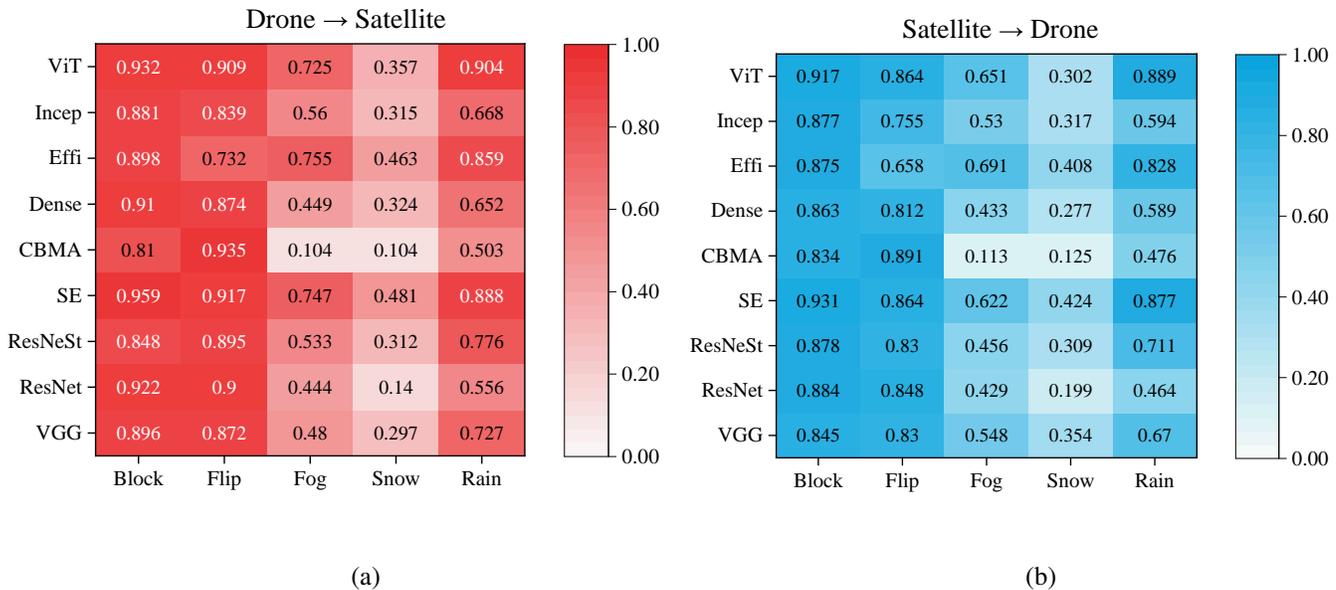


Fig. 10. Heatmaps of the robustness of different backbone networks under unfavorable factors. (a):Drone → Satellite. Darker red indicates that the model was less affected and showed better robustness. (b):Satellite → Drone. Darker blue indicates that the model was less affected and showed better robustness.

TABLE IV  
THE NUMBER OF PARAMETERS OF ALL MODELS AND THEIR INFERENCE SPEEDS WITH BENCHMARK

Backbone	Params(M)	Drone → Satellite	Satellite → Drone
VGG16-bn	272.86	1.18	1.17
ResNet-50	49.24	1.00	1.00
SE-ResNet-50	54.30	1.02	1.02
ResNeSt-50	53.09	1.02	1.00
CBAM-ResNet-50	59.30	1.04	1.02
DenseNet-201	35.73	1.05	1.02
EfficientNetv1-b4	37.06	1.01	1.00
Inceptionv4	83.98	1.03	1.01
ViT-base	172.20	2.45	2.48

may be observed from Figure 9, the accuracy of ViT decreased by only 23.13% and 14.28% of accuracy on two tasks from 300m to 150m, respectively. This shows strong robustness when the size of the target scene changes. We argue that the self-attention mechanism of Transformer architecture helps ViT ignore the redundant information at low heights.

**Robustness to uncertainties.** Drones commonly face many negative factors in the outdoor environment. We simulate these situations with image augmentation techniques. Equation (4) is used to evaluate the influence of those factors. Figure 10 shows the ability of different models to cope with uncertainties. The vertical axis shows that ViT and SE-ResNet exhibit better resistance to uncertainties than other models. The horizontal axis shows that snow was the most difficult factor to overcome. Extracting invariant features from images in wintertime scenes is very challenging; no model is able to reach even the original 50% AP value under this condition.

**Inference speed.** In the inference phase, inference speed is a vital evaluation metric, and it also directly determines whether the model can be put into practical application. Therefore, we evaluate the inference speed of different models under two tasks, as may be observed from Table IV. We take the inference time of ResNet as the baseline time: 1.00.

We find that ViT spent the most time on inference, and Task1 and Task2 were 2.45 and 2.48 times the base time, respectively. Although it has fewer parameters than VGG, the computational complexity of ViT is higher.

### C. Multiple Queries

**We consider whether multi-angle feature fusion improves the efficiency of matching.** In previous matching experiments, a single drone-view image was used as a query for Drone → Satellite. Each scene in the SUES-200 dataset provides a full 360-degree view of the drone view image, which provides complete and comprehensive information about the target scene. Therefore, if a single query cannot describe the target scene, we can introduce multiple angles of drone view images as queries. as may be observed from Table V, we set the multiple-query image to 50, 25, 10, 5, 1. The experimental results show that the multiple queries contain more images, and the Recall@K and AP of the matching are enhanced accordingly. When the average features of 50 images are used as queries, the accuracy of Recall@1 is generally improved by 15%, compared with the single query.

TABLE V  
THE MATCHING ACCURACY (%) OF MULTIPLE QUERIES BASED ON THE BASELINE. 50,25,10,5,1 DENOTE MULTIPLE-QUERY IMAGE SETTING

Query	Drone → Satellite							
	150m		200m		250m		300m	
	Recall@1	AP	Recall@1	AP	Recall@1	AP	Recall@1	AP
50	75.00	78.99	75.00	78.50	78.75	82.27	85.00	87.41
25	71.25	75.81	73.13	77.00	76.25	80.10	85.63	87.58
10	67.75	72.68	69.00	73.50	76.25	79.95	84.25	86.57
5	65.37	70.49	68.25	72.72	75.50	79.22	82.00	84.75
1	59.32	64.93	62.30	67.24	71.35	75.49	77.17	80.67

TABLE VI  
TEST RESULTS OF TRANSFER LEARNING MODELS AND PRE-TRAINED WEIGHTS ON SUES-200.

Training Set	Drone → Satellite							
	150m		200m		250m		300m	
	Recall@1	AP	Recall@1	AP	Recall@1	AP	Recall@1	AP
ImageNet	13.20	17.83	16.70	22.15	13.55	17.96	14.27	18.84
University-1652	54.90	61.11	63.55	68.82	68.53	73.20	72.00	76.29
SUES-200(from scratch)	14.43	19.03	18.25	23.25	21.22	26.36	24.30	29.96
SUES-200(ImageNet pre-trained)	59.32	64.93	62.30	67.24	71.35	75.49	77.17	80.67
SUES-200(U1652 pre-trained)	71.67	75.55	75.57	78.97	79.97	82.50	81.42	84.11
Training Set	Satellite → Drone							
	150m		200m		250m		300m	
	Recall@1	AP	Recall@1	AP	Recall@1	AP	Recall@1	AP
ImageNet	16.25	9.85	7.50	6.38	18.75	11.96	26.25	16.00
University-1652	61.25	48.08	75.00	60.24	77.50	66.51	75.00	70.29
SUES-200(from scratch)	17.50	11.62	30.00	18.56	35.00	22.13	47.50	29.46
SUES-200(ImageNet pre-trained)	82.50	58.95	85.00	62.56	88.75	69.96	96.25	84.16
SUES-200(U1652 pre-trained)	85.00	71.36	86.25	75.96	88.75	79.54	92.50	84.89

#### D. Transfer Learning

We consider whether previous datasets can help a model learn features at different heights, and whether pre-trained weights exhibit an impact on training process. We test whether the models obtained from training on the ImageNet dataset, as well as the University-1652 dataset, can extract discriminative features at different heights. We train a model from scratch on SUES-200 with ImageNet and University-1652 as pre-trained weights. The backbone networks in the above models are ViT. As shown in Table VI, the University-1652-based transfer learning model achieves surprising results compared to ImageNet, which validates that University-1652 can be applied to real scenes. But University-1652’s ability at different heights is still limited because the dataset does not distinguish the effects of different heights. Further, we find that the model trained from scratch is much less capable of extracting features than the model trained based on ImageNet. Another interesting finding is that the initial training process of the model based on the pre-trained weights of University-1652 performed better than the one based on ImageNet, which also shows that the initialization weights of the model are significantly important.

#### E. Other Baseline Models on SUES-200

We evaluate the performance of the classical cross-view matching model on the SUES-200 dataset. Some previous works [17], [19] have designed deep neural networks that achieved excellent performance on different cross-view matching datasets. We select LCM [17] and LPN [19] and migrate their backbone network designs into our pipeline for

training. The experimental results are shown in Table VII. Due to the feature partitioning strategy presented by the LPN for extracting semantic information, the strategy is able to extract global features of the image instead of focusing on the center of the image alone. LPN achieves competitive performance on SUES-200, especially in Task 1.

## VI. ABLATION STUDY

### A. Effect of Feature Dimensions

We also consider how different feature dimensions affect the model. The dimensionality of features extracted from the drone and satellite images in the SUES-200 dataset was 512, as shown in Figure 7, FC1. Therefore, in the ablation learning phase, we reset the feature dimension to 256 and 1024, keeping all other conditions constant. As shown in Table VIII, when we set the dimension to 256, the Recall@1 and AP accuracy both decreased. When we set the dimension to 1024, the performance is better than with 512 dimensions in some metrics. However, 512 is still the overall optimal size and we thus use this dimensionality in the baseline model.

### B. Effects of sharing weights

We consider whether sharing weights could help the model learn better features. With increasing height, drone and satellite images become more and more similar. Hence, the question arises as to whether the model learning efficiency of the model can be improved by sharing the weights of the backbone. We test the effects of sharing model weights on the final test results in the baseline model. Figure 11 shows that the evaluation metrics of both tasks show significant decreases

TABLE VII  
TEST PERFORMANCES OF LCM AND LPN ON SUES-200

Methods	Drone → Satellite							
	150m		200m		250m		300m	
	Recall@1	AP	Recall@1	AP	Recall@1	AP	Recall@1	AP
SUES-200 baseline	59.32	64.93	62.30	67.24	71.35	75.49	77.17	80.67
LCM [17]	43.42	49.65	49.42	55.91	54.47	60.31	60.43	65.78
LPN(block=4) [19]	61.58	67.23	70.85	75.96	80.38	83.80	81.47	84.53
Methods	Satellite → Drone							
	150m		200m		250m		300m	
	Recall@1	AP	Recall@1	AP	Recall@1	AP	Recall@1	AP
SUES-200 baseline	82.50	58.95	85.00	62.56	88.75	69.96	96.25	84.16
LCM [17]	57.50	38.11	68.75	49.19	72.50	47.94	75.00	59.36
LPN(block=4) [19]	83.75	66.78	88.75	75.01	92.50	81.34	92.50	85.72

TABLE VIII  
ABLATION STUDY OF DIFFERENT FEATURE DIMENSIONS ON THE SUES-200 DATASET.

Feature Dimension	Drone → Satellite								
	150m		200m		250m		300m		Height Robustness
	Recall@1	AP	Recall@1	AP	Recall@1	AP	Recall@1	AP	
256	48.15	54.95	57.80	63.14	65.13	70.26	70.15	75.03	31.36%
512	59.32	64.93	62.30	67.24	71.35	75.49	77.17	80.67	23.13%
1024	51.13	56.71	64.35	69.09	72.25	76.34	76.63	80.24	33.28%
Feature Dimension	Satellite → Drone								
	150m		200m		250m		300m		Height Robustness
	Recall@1	AP	Recall@1	AP	Recall@1	AP	Recall@1	AP	
256	63.75	48.86	77.50	60.81	83.75	69.23	86.25	71.54	26.09%
512	82.50	58.95	85.00	62.56	88.75	69.96	96.25	84.16	14.29%
1024	72.50	52.75	83.75	67.52	88.75	76.77	90.00	76.46	19.44%

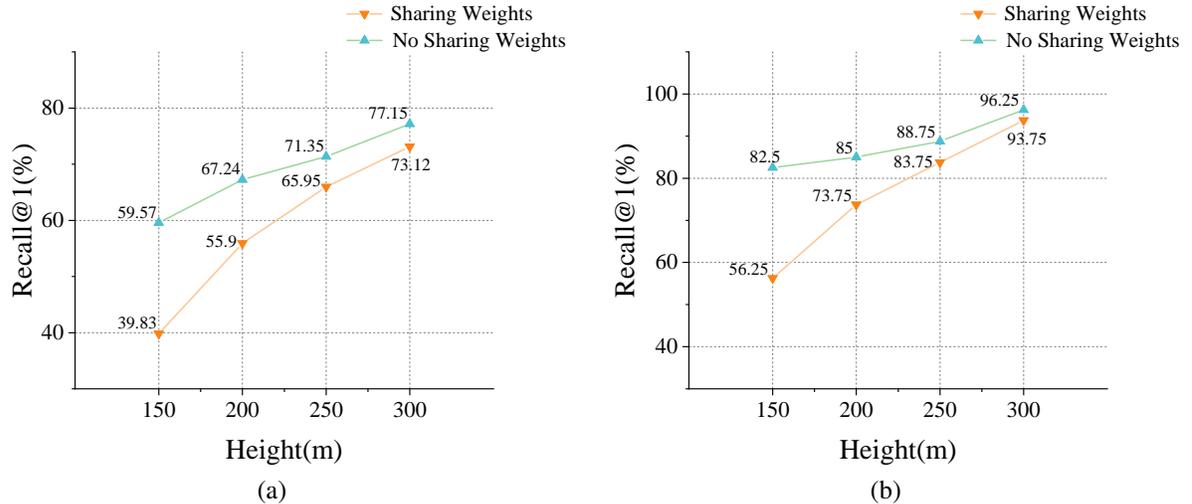


Fig. 11. The accuracy of Recall@1 without sharing weights is always higher than that of Recall@1 with sharing weights, but the gap decreases as the height rises. (a) Drone → Satellite (b) Satellite → Drone

when the sharing weights are available. Still, the difference values between the shared and unshared weights decreased with height. Images collected with increasing height are more similar to satellite images. One possible explanation is that sharing weights can help the model extract more efficient features in similar image pairs.

### C. Effects of different loss function

We also consider whether other loss functions would affect the learning performance of the model. The most common loss functions in previous studies of matching re-

trieval tasks are contrastive loss [43] and triplet loss [44], and these loss functions achieve good performance in other works, such as ReID. To verify the feasibility of these loss functions on our baseline, we strictly hold the backbone network and other parameters constant during the experiments. From Table IX, it may be observed that each of these three loss functions shows advantages and disadvantages in terms of Recall@K and AP. However, when evaluating robustness to height, the accuracy of cross-entropy loss fall the least accuracy from 300m to 150m.

TABLE IX  
ABLATION STUDY OF DIFFERENT LOSS TERMS

Loss	Drone → Satellite								Height Robustness
	150m		200m		250m		300m		
	Recall@1	AP	Recall@1	AP	Recall@1	AP	Recall@1	AP	
CrossEntropy [43]	59.32	64.93	62.30	67.24	71.35	75.49	77.17	80.67	23.13%
Contrastive [44]	54.03	59.85	59.82	63.12	67.28	71.58	71.60	75.79	24.53%
Triplet(margin=0.3)	51.34	57.36	63.57	68.49	68.62	73.13	71.72	75.69	28.41%

Loss	Satellite → Drone								Robustness
	150m		200m		250m		300m		
	Recall@1	AP	Recall@1	AP	Recall@1	AP	Recall@1	AP	
CrossEntropy [43]	82.50	58.95	85.00	62.56	88.75	69.96	96.25	84.16	14.29%
Contrastive [44]	75.00	58.06	81.25	60.31	86.25	70.97	91.25	73.14	17.81%
Triplet(margin=0.3)	75.00	56.95	81.25	59.52	85.00	68.09	87.50	77.08	14.29%

TABLE X  
ABLATION STUDY OF DISTANCE MEASUREMENT ALGORITHMS ON SUES-200

Distance	Drone → Satellite							
	150m		200m		250m		300m	
	Recall@1	AP	Recall@1	AP	Recall@1	AP	Recall@1	AP
Euclidean	59.70	65.24	62.17	67.13	71.30	75.46	77.28	80.75
Manhattan	57.30	62.98	61.83	66.78	69.10	73.61	75.52	79.38
Cosine	59.32	64.93	62.30	67.24	71.35	75.49	77.17	80.67

Distance	Satellite → Drone							
	150m		200m		250m		300m	
	Recall@1	AP	Recall@1	AP	Recall@1	AP	Recall@1	AP
Euclidean	81.25	58.93	85.00	62.55	90.00	69.98	95.00	84.14
Manhattan	76.25	56.80	85.00	59.69	88.75	68.29	92.50	81.98
Cosine	82.50	58.95	85.00	62.56	88.75	69.96	96.25	84.16

#### D. Effects of Distance Measurement Algorithm

We consider whether different distance measurement algorithms affect the matching results. In cross-view matching, several measurement algorithms exist, such as euclidean distance, manhattan distance, and cosine distance. We apply cosine distance in our baseline model due to its good performance in image retrieval tasks [45], [46]. How do other distance measurement algorithms perform on SUES-200? As shown in Table X, Manhattan distance has the worst performance compared with other distance measurement algorithms. Euclidean distance achieves comparable results to the cosine distance.

#### E. Effects of Distractors in Gallery

We also consider whether distractors affect the matching process. In the test stage, we add training data to the gallery set as distractors. To compare the performance of the model without distractors, we remove training data from the gallery set. It can be seen in Table XI, the model's performance improves significantly under the gallery set without distractors. The absence of distractors made it easier for the model to find the correct match. Therefore, we believe a gallery set with more data allows for a more comprehensive evaluation of the model.

#### F. Effects of adding the losses

We consider whether adding the losses would affect model. To verify the effect of adding the losses from two branches, we divide the network into two parts in the training

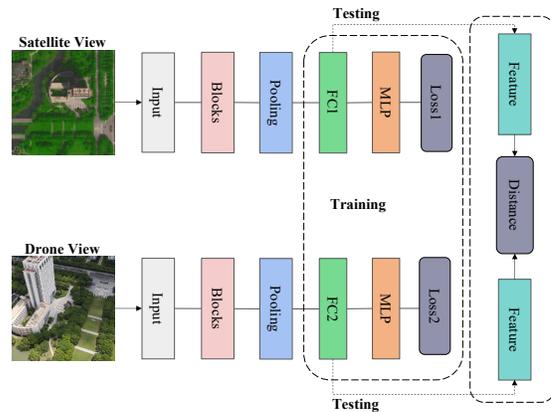


Fig. 12. The structure of the model for calculating the loss separately

stage. As shown in Figure 12, we train two SE-ResNet50 models at 150m height. Their loss functions are no longer added together, and their models were optimized independently. Experimental results are shown in Table XII; the performance of both tasks considerably declined. Why do the results exceed our expectations? How is the current network model different from the previous one? The key is the first FC layer (See FC1 in Figure 7), which shares weights with two branches in the original network. If we also share FC layer weights when testing the current network, the experimental results are presented in Table XII. We initialize FC1 and FC2 with the weight of either view in the current network (Figure 12), the result is another huge improvement compared to the

TABLE XI  
ABLATION STUDY OF DISTRACTORS ON SUES-200

Gallery set	Drone → Satellite							
	150m		200m		250m		300m	
	Recall@1	AP	Recall@1	AP	Recall@1	AP	Recall@1	AP
No Distractors	72.03	76.12	73.18	77.10	80.48	83.31	83.17	85.90
Distractors	59.32	64.93	62.30	67.24	71.35	75.49	77.17	80.67

Gallery set	Satellite → Drone							
	150m		200m		250m		300m	
	Recall@1	AP	Recall@1	AP	Recall@1	AP	Recall@1	AP
No Distractors	88.75	71.65	90.00	70.93	91.25	76.83	95.00	88.29
Distractors	82.50	58.95	85.00	62.56	88.75	69.96	96.25	84.16

TABLE XII  
EFFECT OF ADDING THE LOSSES

FC layer weight	Drone → Satellite			
	Recall@1	Recall@5	Recall@10	AP
Adding losses jointly optimize	56.01	80.12	91.18	62.20
Separated losses independently optimize	4.65	12.00	20.65	7.62
Satellite init	52.97	78.72	89.10	58.87
Drone init	49.20	76.40	87.55	55.42

FC layer weight	Satellite → Drone			
	Recall@1	Recall@5	Recall@10	AP
Adding losses jointly optimize	75.00	86.12	89.58	55.51
Separated losses independently optimize	5.00	7.50	8.75	6.17
Satellite init	62.50	70.00	71.25	53.87
Drone init	50.00	57.50	63.50	44.57

TABLE XIII  
EFFECT OF DIFFERENT ENSEMBLE STRATEGIES

Method	Drone → Satellite							
	150m		200m		250m		300m	
	Recall@1	AP	Recall@1	AP	Recall@1	AP	Recall@1	AP
Average	66.25	71.61	77.50	80.84	80.00	83.84	82.50	85.49
Max pooling	58.74	65.19	72.50	76.19	66.25	72.37	77.50	80.92
Voting	55.00	62.46	71.25	75.15	76.25	79.95	84.25	86.57

former one.

Finally, we summarize the potential advantages of adding losses. Firstly, we believe the two branches constrain each other to jointly optimize the FC Layer, resulting in the FC layer being able to extract features available to both views. Secondly, a complete network system is more suitable for adjusting the hyperparameters in the training stage.

### G. Effects of different ensemble strategies

**We also consider whether different ensemble strategies would affect the multiple queries.** Apart from the numerical average, we also try other ensemble strategies. We employ max pooling and voting to fuse features. Max pooling is a common fusion method. We present it to extract the most remarkable part of the feature maps. The voting method is to select the one with the most occurrences among the predicted labels as the prediction. We apply SE-ResNet-50 as the backbone to evaluate those strategies. The number of query images is 50. It can be seen in Table XIII, we observe that the numerical average arrives better performance than Max pooling and voting on both tasks. We also notice that as the height rises and more predictions become correct, the performance of the voting method also rises.

## VII. VISUALIZATION

Figure 13 shows the visualization results of the baseline model under Rank 5 at different heights and two tasks. It may be observed that ViT (baseline) was able to retrieve the correct result in some very similar scenes. Furthermore, we also visualize the heat maps generated by different models on SUES-200. Figure 14 compares the results of ResNet, Dense, and ViT on the drone view and satellite view. We observe that CNN-based models focused attention on the main target. ViT was also able to notice other things in the background around the target scene. It even had the capability to accurately depict the shape of the target.

## VIII. DISCUSSION

In this study, we find that height exhibited a significant effect on cross-view matching. At the heights of 150m and 200m, the drone footage was more influenced by the surrounding environment and the camera pose. The size of the target scene size leads to drone images being very different from the satellite view images. Hence, accuracy at low heights is relatively low. However, with increasing height, the drone is less influenced by redundant information. The accuracy of matching gradually increases. At the same time, we consider

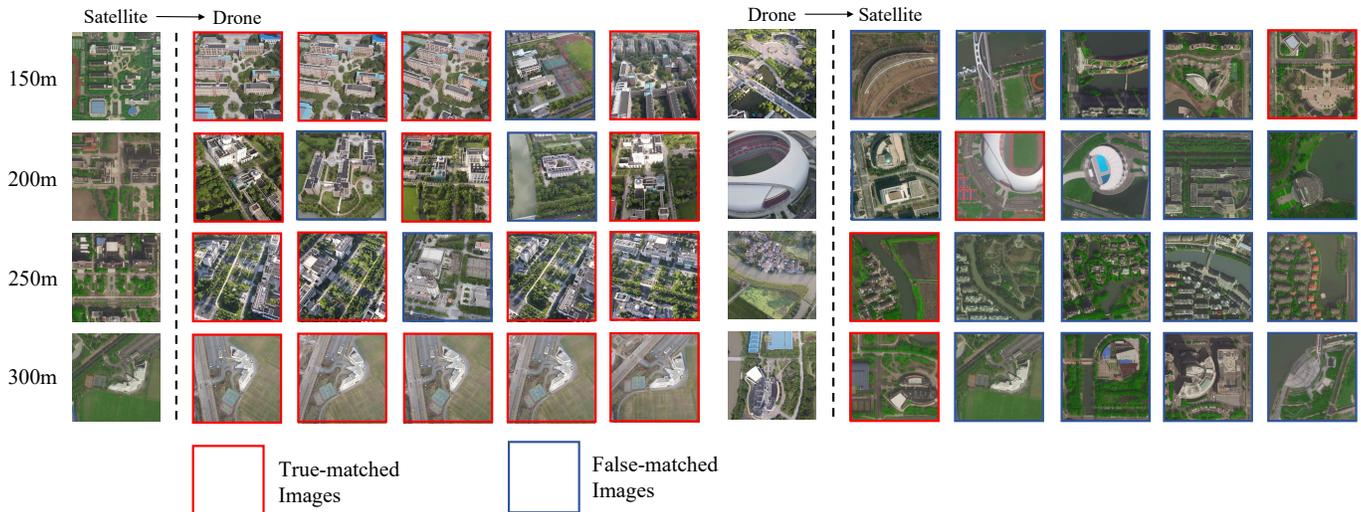


Fig. 13. Qualitative image retrieval results. Top-5 retrieval results of drone view target localization on SUES-200. Top-5 retrieval results of drone navigation on SUES-200.

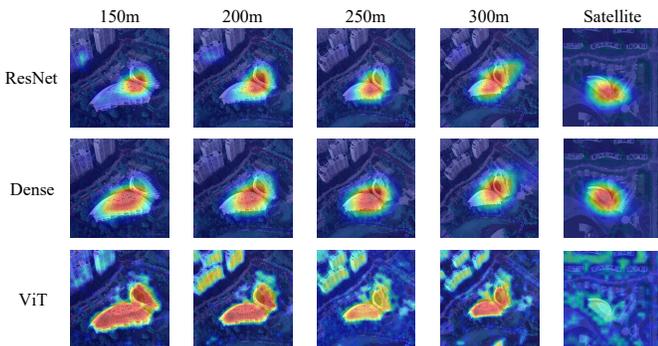


Fig. 14. Visualization of heatmaps. Heatmaps are generated by ResNet, Dense and Baseline based model

that the bottleneck of previous research on cross-view matching studies lies in the lack of a suitable feature extractor. As shown in Table I, most methods are based on the same feature extractor, ResNet.

To test the performance of these feature extractors in a complete way, we also evaluate their performance in the other three aspects through the pipeline. The results show that the ViT-based model has better robustness at different heights and uncertainties. However, it is very challenging for most approaches to extract invariant features in wintertime. Furthermore, the inference speed of the ViT-based model still needs to be improved. Another limitation is that SUES-200 still suffers from a small number of samples and limited viewpoints of the same location.

## IX. CONCLUSION

In this study, we have investigated the problem of image matching across drone and satellite views at different heights. We have proposed a multi-height, multi-scene benchmark called SUES-200, which contains images collected from aerial drones and satellite images for 200 locations. We have also

presented three metrics with a pipeline to comprehensively evaluate the model's ability. 1) robustness at different heights; 2) robustness to uncertainties; 3) inference speed; The results of our experiments have shown that the accuracy and precision of matching increase as the drone's height increases. After evaluating different feature extractors, we provide the model with the best overall performance as the baseline model of SUES-200. For Drone  $\rightarrow$  Satellite, baseline achieves 59.32, 62.30, 71.35, and 77.17 Recall@1 accuracy at heights of 150m, 200m, 250m, and 300m, respectively. For Satellite  $\rightarrow$  Drone, baseline achieves 82.50, 85.00, 88.75, and 96.25 Recall@1 accuracy in 150m, 200m, 250m, and 300m, respectively. We also observe that appropriate pre-trained weights and multiple queries can benefit the model to achieve even better performance, which provides an approach to improve matching efficiency further.

In the future, the main issue to be considered is how to filter out the invalid redundant information at low heights. We also plan to develop a network to adapt to different flight conditions, when the cross-view matching system faces uncertainties. Moreover, the development of a lightweight Transformer architecture for cross-view matching would also be very beneficial to the application of this technology. The data of SUES-200 will also be extended in future research.

## REFERENCES

- [1] T. Y. Lin, S. Belongie, and J. Hays, "Cross-view image geolocation," *Springer International Publishing*, 2016.
- [2] T.-Y. Lin, Y. Cui, S. Belongie, and J. Hays, "Learning deep representations for ground-to-aerial geolocation," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 5007–5015.
- [3] M. Zhai, Z. Bessinger, S. Workman, and N. Jacobs, "Predicting ground-level scene layout from aerial imagery," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 867–875.
- [4] Y. Tian, C. Chen, and M. Shah, "Cross-view image matching for geo-localization in urban environments," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3608–3616.

- [5] Y. Shi, L. Liu, X. Yu, and H. Li, "Spatial-aware feature aggregation for image based cross-view geo-localization," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [6] Y. Shi, X. Yu, L. Liu, T. Zhang, and H. Li, "Optimal feature transport for cross-view image geo-localization," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 11 990–11 997.
- [7] Y. Shi, X. Yu, D. Campbell, and H. Li, "Where am i looking at? joint location and orientation estimation by cross-view matching," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 4064–4072.
- [8] S. Workman and N. Jacobs, "On the location dependence of convolutional neural network features," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2015, pp. 70–78.
- [9] L. Liu and H. Li, "Lending orientation to neural networks for cross-view geo-localization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5624–5633.
- [10] L. Wang, J. Li, B. Huang, J. Chen, X. Li, J. Wang, and T. Xu, "Auto-perceiving correlation filter for uav tracking," *IEEE Transactions on Circuits and Systems for Video Technology*, 2022.
- [11] C. Zhan, H. Hu, X. Sui, Z. Liu, J. Wang, and H. Wang, "Joint resource allocation and 3d aerial trajectory design for video streaming in uav communication systems," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 8, pp. 3227–3241, 2020.
- [12] Y. Sun, B. Cao, P. Zhu, and Q. Hu, "Drone-based rgb-infrared cross-modality vehicle detection via uncertainty-aware learning," *IEEE Transactions on Circuits and Systems for Video Technology*, pp. 1–1, 2022.
- [13] C. ZHAO, Y. ZHOU, Z. LIN, J. HU, and Q. PAN, "Review of scene matching visual navigation for unmanned aerial vehicles," *SCIENTIA SINICA Informationis*, vol. 49, no. 5, pp. 507–519, 2019.
- [14] X. Zhuo, T. Koch, F. Kurz, F. Fraundorfer, and P. Reinartz, "Automatic uav image geo-registration by matching uav images to georeferenced image data," *Remote Sensing*, vol. 9, no. 4, p. 376, 2017.
- [15] D. L. Krishnan, K. Kannan, R. Muthaiah, and M. R. Nalluri, "Evaluation of metrics and a dynamic thresholding strategy for high precision single sensor scene matching applications," *Multimedia Tools and Applications*, vol. 80, no. 12, pp. 18 803–18 820, 2021.
- [16] Z. Zheng, Y. Wei, and Y. Yang, "University-1652: A multi-view multi-source benchmark for drone-based geo-localization," in *Proceedings of the 28th ACM international conference on Multimedia*, 2020, pp. 1395–1403.
- [17] L. Ding, J. Zhou, L. Meng, and Z. Long, "A practical cross-view image matching method between uav and satellite for uav-based geo-localization," *Remote Sensing*, vol. 13, no. 1, p. 47, 2021.
- [18] J. Zhuang, M. Dai, X. Chen, and E. Zheng, "A faster and more effective cross-view matching method of uav and satellite images for uav geo-localization," *Remote Sensing*, vol. 13, no. 19, p. 3979, 2021.
- [19] T. Wang, Z. Zheng, C. Yan, J. Zhang, Y. Sun, B. Zhenga, and Y. Yang, "Each part matters: Local patterns facilitate cross-view geo-localization," *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.
- [20] X. Tian, J. Shao, D. Ouyang, and H. T. Shen, "Uav-satellite view synthesis for cross-view geo-localization," *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.
- [21] M. Dai, J. Hu, J. Zhuang, and E. Zheng, "A transformer-based feature segmentation and region alignment method for uav-view geo-localization," *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.
- [22] H. O. Song, Y. Xiang, S. Jegelka, and S. Savarese, "Deep metric learning via lifted structured feature embedding," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 4004–4012.
- [23] S. Workman, R. Souvenir, and N. Jacobs, "Wide-area image geo-localization with aerial reference imagery," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 3961–3969.
- [24] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International journal of computer vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [25] H. Bay, T. Tuytelaars, and L. V. Gool, "Surf: Speeded up robust features," in *European conference on computer vision*. Springer, 2006, pp. 404–417.
- [26] J. Revaud, P. Weinzaepfel, C. De Souza, N. Pion, G. Csurka, Y. Cabon, and M. Humenberger, "R2d2: repeatable and reliable detector and descriptor," *arXiv preprint arXiv:1906.06195*, 2019.
- [27] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "Orb: An efficient alternative to sift or surf," in *2011 International conference on computer vision*. Ieee, 2011, pp. 2564–2571.
- [28] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [29] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.
- [30] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [31] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.
- [32] N. N. Vo and J. Hays, "Localizing and orienting street views using overhead imagery," in *European conference on computer vision*. Springer, 2016, pp. 494–509.
- [33] T.-Y. Lin, Y. Cui, S. Belongie, and J. Hays, "Learning deep representations for ground-to-aerial geolocalization," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 5007–5015.
- [34] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [35] H. Zhang, C. Wu, Z. Zhang, Y. Zhu, H. Lin, Z. Zhang, Y. Sun, T. He, J. Mueller, R. Manmatha *et al.*, "Resnet: Split-attention networks," *arXiv preprint arXiv:2004.08955*, 2020.
- [36] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "Cbam: Convolutional block attention module," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 3–19.
- [37] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [38] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International Conference on Machine Learning*. PMLR, 2019, pp. 6105–6114.
- [39] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *Thirty-first AAAI conference on artificial intelligence*, 2017.
- [40] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [41] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1026–1034.
- [42] A. B. Jung, K. Wada, J. Crall, S. Tanaka, J. Graving, C. Reinders, S. Yadav, J. Banerjee, G. Vecsei, A. Kraft, Z. Rui, J. Borovec, C. Vallentin, S. Zhydenko, K. Pfeiffer, B. Cook, I. Fernández, F.-M. De Rainville, C.-H. Weng, A. Ayala-Acevedo, R. Meudec, M. Laporte *et al.*, "imgaug," <https://github.com/aleju/imgaug>, 2020, online; accessed 01-Feb-2020.
- [43] R. Hadsell, S. Chopra, and Y. LeCun, "Dimensionality reduction by learning an invariant mapping," in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, vol. 2. IEEE, 2006, pp. 1735–1742.
- [44] E. Hoffer and N. Ailon, "Deep metric learning using triplet network," in *International workshop on similarity-based pattern recognition*. Springer, 2015, pp. 84–92.
- [45] X. Zhu, Z. Luo, P. Fu, and X. Ji, "Voc-reid: Vehicle re-identification based on vehicle-orientation-camera," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 602–603.
- [46] H. Luo, Y. Gu, X. Liao, S. Lai, and W. 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.



**Runzhe Zhu** received the B.S. degree in Zhejiang Shuren University from Zhejiang, Hangzhou, China, in 2020. He is currently an M.S. student with the department of electrical and electronic engineering of Shanghai University of Engineering Science, Shanghai, China. His research interests include visual geolocalization and cross-view matching.



**Wenbo Hu** received the Ph.D. degree in Geography from Université Grenoble Alpes, Grenoble, France, in 2019. He is currently the post-doctor in School of Communication and Information Engineering, Shanghai University. From 2019 to 2020, he was a postdoctoral researcher with the Laboratoire PACTE, UMR 5194 CNRS, France. His research interest include behavioral geography, spatial modeling and public policing based on big data, deep learning and machine learning.



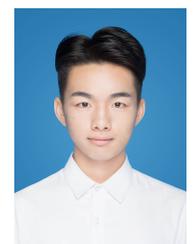
**Ling Yin** received the B.S. degree in software engineering from East China Normal University, China, in 2008, and the Ph.D. degree in Computer technology from East China Normal University, China, in 2016. She is currently a lecturer with the School of Electronic and Electrical Engineering, Shanghai University of Engineering Science, China. Her research interest includes deep learning, time series analysis, and software engineering with formal methods.



**Mingze Yang** received the B.S. degree in Rolling Stock Engineering from Shanghai Institute of Technology, Shanghai, China, in 2019. He is currently a M.S. student with the department of electrical and electronic engineering of Shanghai University of Engineering Science, Shanghai, China. His research interests include deep learning, target recognition and wireless sensing.



**Fei Wu** received the B.S. degree, the M.S. degree and the Ph.D. degree in Computer Science from National University of Defense Technology in 1990, 1993 and 1998. He was a Post-Doctoral Research with Nankai University, China. He is currently a full professor with the School of Electronic and Electrical Engineering, Shanghai University of Engineering Science, China. His research interests include intelligent information processing, positioning technology and machine learning.



**Yuncheng Yang** received the B.S. degree in mechanical engineering from Henan University of Science and Technology, Luoyang, China, in 2018. He is currently a M.S. student with the department of electrical and electronic engineering of Shanghai University of Engineering Science, Shanghai, China. His research interests include deep learning, indoor positioning and wireless sensing.