

# Cross-View Image Sequence Geo-localization

Xiaohan Zhang<sup>1</sup>, Waqas Sultani<sup>2</sup>, and Safwan Wshah<sup>1</sup>

<sup>1</sup>Department of Computer Science, University of Vermont, USA

<sup>2</sup>Intelligent Machine Lab, Information Technology University, Pakistan

{Xiaohan.Zhang, Safwan.Wshah}@uvm.edu waqas.sultani@itu.edu.pk

## Abstract

*Cross-view geo-localization aims to estimate the GPS location of a query ground-view image by matching it to images from a reference database of geo-tagged aerial images. To address this challenging problem, recent approaches use panoramic ground-view images to increase the range of visibility. Although appealing, panoramic images are not readily available compared to the videos of limited Field-Of-View (FOV) images. In this paper, we present the first cross-view geo-localization method that works on a sequence of limited FOV images. Our model is trained end-to-end to capture the temporal structure that lies within the frames using the attention-based temporal feature aggregation module. To robustly tackle different sequences length and GPS noises during inference, we propose to use a sequential dropout scheme to simulate variant length sequences. To evaluate the proposed approach in realistic settings, we present a new large-scale dataset containing ground-view sequences along with the corresponding aerial-view images. Extensive experiments and comparisons demonstrate the superiority of the proposed approach compared to several competitive baselines.*

## 1. Introduction

Cross-view image geo-localization aims to determine the geospatial location from where an image was taken (also known as the query image) in a database of geo-tagged aerial images (also known as reference images) [40, 30, 19, 43]. Estimating geo-spatial locations from images has many important applications such as autonomous driving [29], robot navigation[4, 17], augmented reality (AR) [9], and unmanned aerial vehicle (UAV) navigation [29].

Despite the huge research efforts that have been done on this problem, image geo-localization remains far from being solved and is considered one of the most challeng-

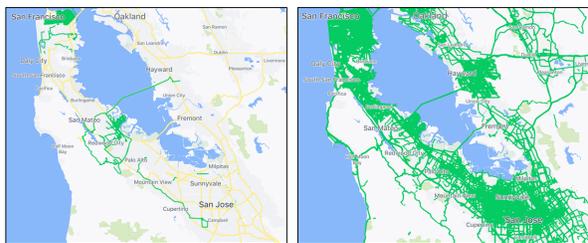


Figure 1: Comparison of the coverage area (green lines) of user uploaded street view images between panoramas (left) and limited FOV images (right) in San Francisco, USA from Mapillary [2].

ing tasks in the computer vision field due to: 1) the drastic appearance differences between the query images and reference images, 2) capturing time gaps between the query image and the reference image results in different illumination conditions, weather, and objects and, 3) differences in resolution at which ground and aerial images are captured.

Recent research in cross-view image geo-localization has shown tremendous progress on large-scale datasets [40, 21, 43], but they heavily rely on panoramic query images [40, 16, 30, 26, 32, 6, 21, 42, 43, 39]. Even though panoramic images provide richer contextual information than normal limited Field-Of-View (FOV) images, in practice, limited FOV images are more common and easier to capture from smartphones, dash cams, and digital single-lens reflex (DSLR) cameras. Fig. 1 shows the comparison of coverage area of users uploaded street view images between panoramas and limited FOV images in San Francisco, USA on Mapillary [2]. Moreover, even map platforms such as Google Street View (GSV) provides panoramas only for a few historic or tourist attraction places for several countries such as China, Qatar, and Pakistan. However, limited FOV street view images are available across 190 countries in the most of the regions as shown on Mapillary [2]. Clearly, the limited FOV images are much more popular than panoramic

images. This applies to all other countries and is more noticeable in developing countries where for the most part, panoramic images are not available at all.

Due to the recent advancement of autonomous vehicles and the Advanced Driving Assistance System (ADAS), frontal street view videos are easily accessible from the dash cams in current vehicles. Instead of using unpopular panoramic images [31, 35, 30], expanding cross-view geo-localization algorithms to work on sequences of images is more practical and more acceptable in real-world scenarios. On the other hand, current cross-view geo-localization approaches [31, 35, 30, 16, 32, 43, 37] deal mainly with a single image for geo-localization and cannot be used directly to capture the temporal structure that lies within a sequence of FOV frames. Thus, it is a natural extension to expand cross-view geo-localization methods on sequences of limited FOV images named cross-view image sequence geo-localization.

In this paper, a new cross-view geo-localization approach is proposed that works on sequences of limited FOV images. Our model is trained end-to-end to capture *temporal feature representations* that lie within the images for better geo-localization. Although our model is trained on fixed-length temporal sequence, it tackles the challenge of variable length sequence during the inference phase through a novel sequential dropout scheme. To the best of our knowledge, we are the first one to propose end-to-end cross-view geo-localization from *sequences* of images. We refer to this task as *cross-view image sequence geo-localization*. Furthermore, to facilitate future research in cross-view geo-localization from sequences, we put forward a new dataset and compare our proposed model with several recent baselines. In summary, our main contributions are as follows:

- 1) We propose a new end-to-end approach, *cross-view image sequence geo-localization*, that geo-localizes a query sequence of limited FOV ground images and its corresponding aerial images.
- 2) We introduce the first large-scale cross-view image sequence geo-localization dataset.
- 3) We propose a novel temporal feature aggregation technique that learns an end-to-end feature representation from a sequence of limited FOV images for sequence geo-localization.
- 4) We propose a new sequence dropout method to predict coherent features on sequences of different lengths. The proposed dropout method helps in regularizing our model and achieve more robust results.

## 2. Related Work

**Cross-view Image Geo-localization:** Before the deep learning era, cross-view image geo-localization methods were based on hand-crafted features [19, 8] such as HoG [10], GIST [24], self-similarity [28], and color his-

tograms. Conventional methods struggled with matching accuracy because of the quality of the features. Due to the resurgence of deep learning in numerous computer vision applications, several deep learning based geo-localization methods [40, 20, 37] have been proposed to extract features from fined-tuned CNN models to improve the cross-view geo-localization accuracy. More recently, Hu *et al.* [16] proposed to aggregate features by NetVLAD [3] layer which achieved significant performance improvements. Shi *et al.* [32] proposed a feature transport module for aligning features from aerial view and street view images. Liu *et al.* [21] explored fusing orientation information into the model which boosted performance. With the development of Generative Adversarial Networks (GANs) [14], Regmi *et al.* [26] proposed a GAN-based cross-view image geo-localization approach using a feature fusion training strategy. Zhu *et al.* [43] recently proposed a new approach (VIGOR) that does not require a one-to-one correspondence between ground images and aerial images. It is also worth mentioning that some methods [30, 31, 35] based on ground-level panorama employ the polar transformation which bridges the domain gap between reference images and query images by prior geometric knowledge. By leveraging this prior geometric property, Shi *et al.* [30] proposed Spatial Aware Feature Aggregation (SAFA) which improves the results on CVUSA [40] and CVACT [21] by a large margin. Similar to [26], Toker *et al.* [35] combined SAFA [30] with a GAN. Their proposed method achieved state-of-the-art results on CVUSA [40] and CVACT [21]. However, to perform the polar transformation, the query image is assumed to be aligned at the center of its reference aerial image which is not always guaranteed in real-world scenarios. The above-mentioned methods rely on panoramic ground-level images. By contrast, our method used more easily available limited FOV images.

We noticed that some previous works [31, 37, 34] studied the cross-view image geo-localization problem using a single limited FOV image as a query. Tian *et al.* [34] proposed a graph-based method that matches the detected buildings in both ground images and aerial images. This method was only applicable in metropolitan areas which contain dense buildings. DBL [37] proposed by Vo *et al* focused on geo-localizing the scene in the image rather than the location of the camera. Dynamic Similarity Matching proposed by Shi *et al.* [31] required polar transformed aerial images as input. Compared to these methods, we assume neither aligned ground-level images nor that our method only works in metropolitan areas. Furthermore, instead of geo-localizing a single limited FOV image, our approach geo-localizes a sequence of limited FOV images.

Recently, Regmi and Shah [27] proposed to geo-localize video sequences in the same-view setting by using a geo-temporal feature learning network and a trajectory smooth-

ing network. On the other hand, in this paper, we incorporate aerial images and ground video sequences to address the problem of cross-view image sequence geo-localization by proposing a transformer-based model. Current cross-view geo-localization approaches can be used for sequential cross-view geo-localization trivially by applying them frame by frame as proposed in [17]. However, we propose an end-to-end approach that automatically processes a whole sequence of images and correlates their features with the corresponding aerial image by building a better feature representation in both temporal and spatial domains. We have compared our results with the best models in the literature that could be applied to our dataset as discussed in the experiments section.

**Transformer/multi-head attention:** Recently, Vaswani *et al.* [36] proposed the transformer module and demonstrated its ability in catching the temporal correlation in time series data. Using the transformer, several works [22, 5, 12] achieved remarkable results in natural language processing tasks. In computer vision, transformers have been used for image classification [13], video segmentation [38], object detection [7], and same-view video geo-localization [27]. In this paper, we combined the transformer with the cross-view image sequence geo-localization to effectively utilize the full range of visibility from the sequential data. Our experiments showed that the transformer can learn to fuse and summarize several features from a sequence of images and predict robust results.

## 3. Dataset

### 3.1. Previous Datasets

Many datasets have been proposed for cross-view image geo-localization [40, 21, 43, 37]. Vo *et al.* [37] proposed a large-scale cross-view geo-localization dataset consisting of more than 1 million pairs of satellite-ground images. The authors collected aerial images from Google Maps and the corresponding ground images from Google Street View from eleven different US cities. Workman *et al.* [40] proposed a Cross-View USA (CVUSA) dataset containing more than 1 million ground-level images across the whole USA. Later, Zhai *et al.* [41] refined the CVUSA dataset by pairing 44,416 aerial-ground images and this has become one of the most popular datasets in this field. In this paper, we refer to this refined version as CVUSA. CVACT [21] followed the same structure as CVUSA and had the same number of training samples as CVUSA but had 10 times more testing pairs. Recently, Zhu *et al.* [43] proposed the VIGOR dataset which is the first non one-to-one correspondent cross-view image geo-localization dataset collected randomly from four major US cities. In order to have systems for practical scenarios in which the queries and reference images pairs are not guaranteed to be always

perfectly aligned, VIGOR defined ‘positive’ and ‘semi-positive’ ground images in one single aerial image. Note that current cross-view geo-localization datasets cannot be easily converted to sequential dataset. To the best of our knowledge, there is no existing dataset that provides *sequential* ground-level images and their corresponding aerial images for cross-view image geo-localization.

### 3.2. Proposed Dataset

Since existing cross-view geo-localization datasets [37, 40, 21, 43] contain only discrete ground images, we collected a new cross-view image sequence geo-localization dataset containing limited FOV images which are much more available and applicable for real-world systems. Table 1 demonstrates the comparison of our proposed dataset with the existing cross-view image geo-localization datasets. Below, we first explain the procedures we followed to collect the ground-level images and then describe the process of capturing aerial imagery.

#### 3.2.1 Ground-Level Imagery

Our data was collected using the Fugro Automatic Road Analyzer (ARAN)<sup>1</sup> which is a road data capturing vehicle capable of collecting different data modalities such as image, LiDAR, and pavement laser. ARAN is also equipped with a GPS and an inertial measurement unit (IMU) sensor for providing precise GPS locations and camera poses. The raw dataset contains over 5000km of urban and suburban roads, and highways in both directions in the state of Vermont, US. In our dataset, we only used the frontal camera images with a resolution of  $1920 \times 1080$ . The distance between each capture point is approximately  $8m$  and the FOV of the camera is around  $120^\circ$ . GPS location and camera heading (compass direction) are also provided for each ground-level image. To represent more real-world scenarios, our dataset contains approximately 70% of images from suburban areas and 30% from urban areas which may be collected from one or two-way driving directions. The ratio of the collected two-way driving direction data is around 30% in which the same street images are captured from both driving directions, for example, north-to-south and south-to-north. The total number of ground images is 118,549 resulting in 38,863 aerial pairs as explained in the following sections. Our dataset covers around 500 kilometers of roads in Vermont. Please refer to the supplementary material for more information.

#### 3.2.2 Sequence Formation

After obtaining the raw ground-level data as described in section 3.2.1, long sequences of raw data should be seg-

<sup>1</sup><https://www.fugro.com/our-services/asset-integrity/roadware/equipment-and-software>

Dataset Comparison	Vo [37]	CVACT [21]	CVUSA [40]	VIGOR [43]	Ours
# of Aerial Images	> 1M	128, 334	44,416	90,618	38,863
# of Ground-level Images	> 1M	128, 334	44, 416	105, 214	118, 549
# of Ground Images per aerial image	1	1	1	~5	~7
Coverage	Urban, Suburb	Urban, Suburb	Urban, Suburb	Urban	Urban, Suburb
Seamless Sampling	No	No	No	Yes	Yes
Sequential Ground-level Images	No	No	No	No	Yes
Orientation	Yes	Yes	Yes	Yes	Yes
Ground-level GPS Location	No	Same	Same	Arbitrary	Arbitrary

Table 1: Comparison between our proposed dataset and other existing cross-view image geo-localization datasets.

mented into several small sequences to be used for cross-view geo-localization. A simple but effective greedy algorithm is employed. Given the raw ground-level image sequence data as  $S = s_0, s_1, \dots, s_N$  where  $N$  is the number of ground-level images that would be segmented into sequence splits. It is not required to keep the same number of images in each split since in real-world scenarios the number of images in a sequence can be variable due to different hardware or software configurations. However, to perform the retrieval task, it is required that the images of any resulting sequence must *lay within one single aerial image*. Our algorithm iterates through each image in  $S$ , denoting the first image as  $s_0$ . After that the distance between  $s_0$  and  $s_t$  is calculated. If the distance is less than a preset threshold value  $\Delta$ , we step to the next image  $s_{t+1}$ . Otherwise,  $[s_0, s_t]$  is the segmented sequence. Then the image in the middle of  $s_0$  and  $s_t$  is set as starting point for the next segment. This process is visualized in Fig. 2. The circles displayed with the same color are segmented in one sequence. If one circle has multiple colors, this circle co-exists in two or more segments. We empirically choose  $\Delta = 50 m$  to ensure that images in any sequence fall within one aerial image with zoom level 20. To make the training procedure consistent and simple, which will be discussed later, we removed 72 sequences which contain less than 7 images. This finally results in 38, 863 sequences with an average of seven images per segment. This sequence formation strategy guarantees that the formed sequence is covered by a single aerial image and does not need to know the length of the raw data. Note that our approach needs seven frames during training. However, we do not have such restrictions while in testing. In real-world scenarios, the distance between each frame may vary and one can simply use techniques such as IMU sensors or visual odometry [23] to estimate the distance between frames.

### 3.2.3 Aerial Imagery

Google Maps Static API [1] is employed to obtain the aerial images for each sequence. Assuming the ground images in a single sequence are on a planar surface, we can determine the geometric center (the arithmetic mean location) of the

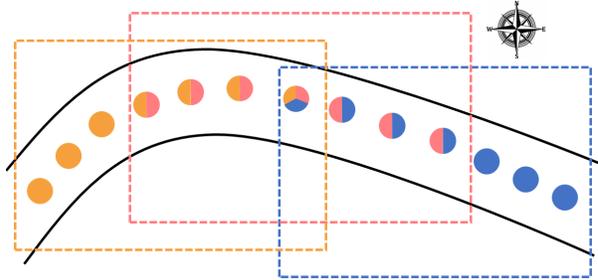


Figure 2: Demonstration of our ground-level images sampling strategy. In this example, three aerial images are captured (yellow, pink, and blue boxes) based on the locations of the ground images (colored circles). Each circle inside these boxes belongs to that aerial image. If one circle has multiple colors, it belongs to multiple sequences.



Figure 3: A sequence sample from our dataset. The aerial image is in the center and the ground images are at the edges. Each orange dot represents the location of one ground image indicated by the blue arrow. The grey arrow of each dot represents the heading direction of the camera.

aerial image for a given sequence. In this way, the aerial image can cover the whole sequence. A random shift at most 5 meters is applied to each aerial image to simulate real-world scenarios. This results a one-to-one correspondence between ground sequences and aerial images. The total number of collected aerial images is 38, 863. Follow-

ing VIGOR [43], each aerial image is captured at the zoom level of 20 with an resolution of  $640 \times 640$ . The ground resolution is approximately 0.114m. A sample pair of ground-aerial images from our dataset is shown in Fig. 3.

## 4. Proposed Methods

### 4.1. Overview

Given a sequence of limited FOV ground images, our goal is to geo-localize an aerial image from a reference database from where this sequence was taken. To achieve this goal, we considered geo-localization as a retrieval task similar to many other previous works [30, 43, 16, 37, 20, 21, 35]. Specifically, we denote extracted features from a geo-tagged aerial image as  $F_{sat}$  and extracted aggregated features from a sequence of ground-level images as  $F_{grd}$ . By evaluating the distance between  $F_{sat}$  and  $F_{grd}$ , we can find the most similar aerial image from a database of aerial images. To extract features from a sequence of ground-level images, we introduce an end-end model to extract the sequential spatio-temporal features. We use VGG16 [33] to extract the spatial features from each image and then pass those features to a novel Temporal Feature Aggregation Module (TFAM) to capture the temporal information. The spatio-temporal features are then aggregated into a single feature for retrieval. Furthermore, to generalize the proposed method for different sequence lengths, a sequential dropout (SD) scheme is implemented. Fig. 4 provides an overview of the proposed approach. In the next sections, we describe TFAM in more detail followed by an introduction to the sequential dropout (SD) scheme in Section 4.3. Finally, in Section 4.4, we describe the training objectives.

### 4.2. Temporal Feature Aggregation Module

To explore the benefits of the sequential images, we introduce TFAM in cross-view sequence geo-localization. TFAM is inspired by the success of the transformers [36] in many computer vision problems [13, 38, 27]. The multi-head self-attention mechanism is the key component that enabled transformers to capture correlations between sequential data elements at any distance. Similar to the transformer [36], TFAM also employs a multi-head self-attention mechanism to capture contextual information from a sequence of images.

Consider a sequence of images as  $\mathcal{P} \in \mathbb{R}^{T \times W \times H \times C}$  where  $T, W, H, C$  are the number of images in a sequence, image width, image height, and image channel respectively. We choose the VGG16 backbone to extract embedding features for each image in the sequence to have a fair comparison with the baseline methods [30, 43]. A feature vector  $F' \in \mathbb{R}^{T \times D}$  is obtained by concatenating each image’s feature along the temporal axis, where  $D$  is the dimension of the output of the backbone feature extractor.

Similar to the original transformer [36], before the multi-head self-attention layer, a sinusoidal positional encoding  $E_{pos} \in \mathbb{R}^{T \times D}$  is added to the extracted feature embeddings  $F'$  to preserve the order of the temporal information which is shown in Equation 1.

$$\tilde{F} = F' + E_{pos}. \quad (1)$$

By feeding the feature embeddings to the Multi-head self-attention layer, each embedding vector is projected into three sub-spaces as  $Q_i = \tilde{F}W_i^Q, K_i = \tilde{F}W_i^K, V_i = \tilde{F}W_i^V$  representing the query, key, value respectively and  $i$  is the index of the head which we will describe later. Note that  $W_i^Q \in \mathbb{R}^{D \times \frac{D}{N_{head}}}, W_i^K \in \mathbb{R}^{D \times \frac{D}{N_{head}}},$  and  $W_i^V \in \mathbb{R}^{D \times \frac{D}{N_{head}}}$  are three projection matrices and  $N_{head}$  is the number of heads. The attention mechanism can be written as follows:

$$head_i = softmax\left(\frac{Q_i K_i^T}{\sqrt{D}}\right) V_i. \quad (2)$$

To fully explore the contextual information in the temporal domain, we use an approach similar to [36] and concatenate the value of the multiple heads and projected into the output space using a projection matrix  $W^O \in \mathbb{R}^{N_{head} D \times D},$

$$F_{aggregated} = Concat(head_1, head_2, \dots, head_{N_{head}}) W^O. \quad (3)$$

By stacking  $N$  TFAM modules, our model can extract more refined feature representations. Finally, the features from the last TFAM,  $F_{aggregated} \in \mathbb{R}^{T \times D}$  has the same input shape as the embedding vector  $F$ . The resulting features are then averaged using an average pooling layer on the temporal axis to obtain a one-dimension vector for the retrieval task as follows:

$$F_{grds} = average\_pool(F_{aggregated}). \quad (4)$$

### 4.3. Adaptive Sequence Length

TFAM introduced in the previous section works well on sequences that have a fixed length  $T$ . However, during inference in real-world settings, it is not always possible to capture exactly  $T$  frames in a sequence due to different hardware or software configurations (e.g. different sampling and capturing rate, signal loss, etc). To adapt TFAM to variant sequence lengths, we propose a sequential dropout (SD) scheme by modifying the TFAM algorithm. During training, a random binary mask  $A \in \mathbb{R}^T$  is generated and fed to each TFAM in the model. For each index  $x$  at  $A_x$ , if  $A_x = 0$ , it means that the feature at index  $x$  in  $\tilde{F}_x$  is omitted. Otherwise, the TFAM operates normally on this feature. By setting  $K_{i,x}$  to a zero vector at index  $x$ , the attention value at head  $i$  of index  $x$  represented as  $Q_i K_{i,x}^T$

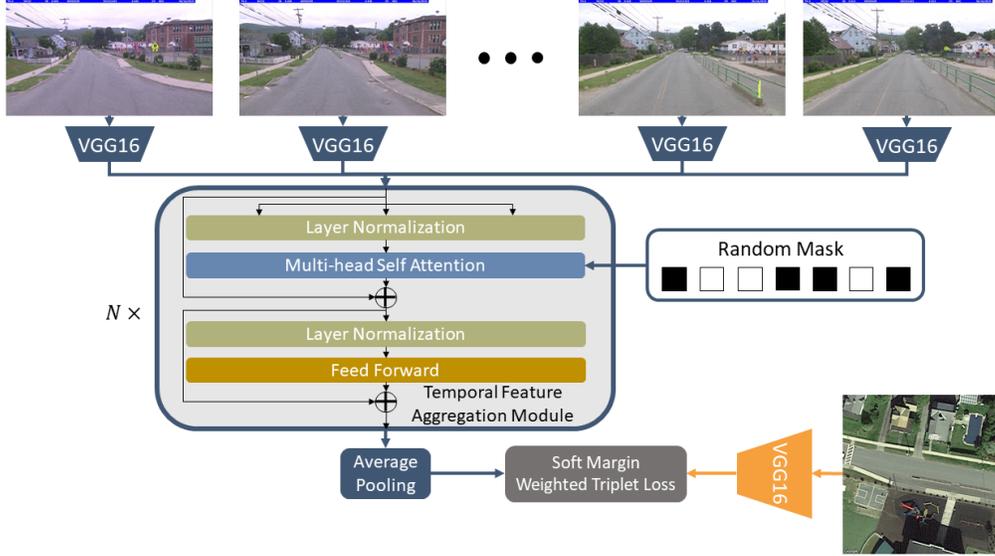


Figure 4: An overview of our proposed method which contains two main parts. The ground features extraction branch (components in dark blue), and the aerial features extraction branch (components in orange). The ground-level features extraction branch takes a sequence of images as input. The aerial features extraction takes the aerial image as input.

also becomes a zero vector. In other words, all the query values would never interact with the key values of this feature at index  $x$ . Consequently, the embedding vector  $\tilde{F}_x$  at index  $x$  of  $\tilde{F}$  is ignored by all other vectors during forward propagation and back-propagation. To generate the random mask  $A$ , we set a maximum number of dropout features  $J$  in which  $J < T$ . In each training mini-batch, we uniformly sampled an integer  $e$  between  $[0, J]$  to represent the number of dropped features in this batch. To control the dropout rate during the training, we initialized all  $A$  values with 1, and randomly set  $e$  elements in  $A$  to 0. The average pooling layer (mentioned in section 4.2) only operates on the index of the temporal dimension of aggregated features  $F_{aggregated}$  where the mask value is 1. Note that to fully exploit the temporal information, our approach does need a fixed-length sequence during training, however, it employs SD to tackle variable-length sequence during testing. In our experiments, we found that this strategy not only helped the TFAM to produce a coherent representation but also regularized the model and achieved much higher performance.

#### 4.4. Training Objective

After extracting the aerial features  $F_{sat}$  and ground-level features  $F'$ ,  $F'$  is further refined using the proposed TFAM as described in 4.2, and the aggregated ground-level features  $F_{grd}$  are obtained. Finally, we deploy a metric learning objective to train the model using weighted soft margin triplet loss [16],

$$\mathcal{L} = \log(1 + e^{\gamma(d_{pos} - d_{neg})}), \quad (5)$$

where  $\gamma$  is a hyperparameter that controls the scale of the loss value.  $d_{neg}$  and  $d_{pos}$  are  $L_2$  distances of unmatched and matched aerial-ground pairs. We employ  $L_2$  normalization on  $F_{sat}$  and  $F_{grd}$  before calculating the distance. The goal of this loss function is to push the matched pairs closer while pushing unmatched pairs further.

## 5. Experiments

**Implementation Details & Dataset:** The proposed method was implemented in PyTorch [25]<sup>2</sup>. We used VGG16 [33] pretrained on ImageNet [11] as backbones of our features extractors. The last two fully connected layers were removed for extracting features. We stacked 6 TFAMs, each with 8 heads, in our model. We adopted our proposed SD scheme with a maximum number of dropout features  $J = 6$  during training. During the testing, the frames can be dropped by setting the values at corresponding locations in  $A$  to 0. Since our proposed method exploit sequence images for training, we cannot evaluate our method on existing cross-view geo-localization datasets. Instead, we chose to benchmark our proposed method on our dataset described in Section 3. The dataset is split into training and testing sets with 31,091 and 7,772 aerial-ground sequence pairs respectively. The training and testing datasets are geographically separated that no overlapping areas between these two datasets. These settings are applied to all the experiments in this section unless specified otherwise.

**Baseline Methods:** We compare our method with

<sup>2</sup>Codes available at <https://gitlab.com/vail-uvn/seqgeo>

	R@1	R@5	R@10	R@1%
VIGOR [43]	0.54%	2.52%	4.48%	18.55%
SAFA <sup>†</sup> [30]	0.68%	2.92%	5.06%	21.81%
SAFA [30]	0.63%	2.83%	5.03%	21.51%
Ours w/o SD	1.39%	<b>6.50%</b>	<b>10.45%</b>	32.42%
Ours w/ SD	<b>1.80%</b>	6.45%	10.36%	<b>34.38%</b>

Table 2: Comparison between our methods with SD and without SD, SAFA and VIGOR methods. † indicates testing on single center ground image as query.

SAFA [30] and VIGOR [43] on our dataset. We chose SAFA [30] as it achieved very competitive results on both CVUSA [40] and CVACT [21] datasets. VIGOR [43] also achieved outstanding performance on their proposed dataset in a one-to-many retrieval approach. To adopt SAFA [30] on our dataset, we trained SAFA [30] on center ground-level images with their corresponding aerial images on our proposed dataset with configurations reported in its original paper. Thus, to make the comparison fair, we initialized SAFA with pre-trained weights on the CVUSA dataset. To be noticed, we did not apply the polar transformation in SAFA [30] to keep a fair comparison with other methods. To train VIGOR, we set the center ground-level image as a ‘positive’ sample and the others are ‘semi-positive’ samples as defined in their original paper. To enable SAFA and VIGOR to work on sequences of images, we feed each ground-level image in the sequence separately and average the final feature vectors for all the images.

**Evaluation Metrics:** Similar to the previous works [30, 16, 21, 43], we use the recall accuracy at top-K (R@K) for evaluating the performance. Given a query sequence, if the ground truth aerial image ranks in the first  $K$  most similar aerial images, it is considered to be a ‘correct’ query.

### 5.1. Quantitative comparison

Our main results are reported in Table 2. SAFA(center) indicates that the SAFA model was tested on the center ground-level image only. SAFA(sequence) means that the SAFA model was tested on the whole sequence of ground-level images after averaging the features. We also provide the results from our method without SD. It can be seen that our method outperforms the baseline methods by a large margin. We also observe that our method performs much better with SD in both top-1 and top-1% recall. The model without SD is slightly better than the model with SD on top-5 and top-10 recalls, but it is a slight margin as shown in the recall vs top-K graph in Fig. 5. Two randomly selected aerial-ground sequence pairs predicted by our model training with SD from our test set are visualized in Fig. 6. In the top two rows of Fig. 6, the ground truth image was successfully predicted as the most similar one. It is worth

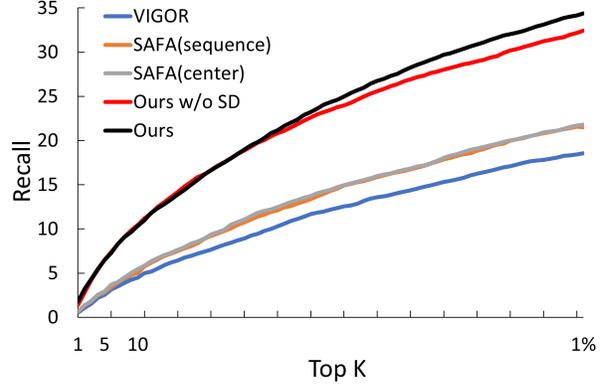


Figure 5: Recall rates of our methods vs baseline methods. The results demonstrate that both methods trained with SD and without SD outperform the baseline methods.



Figure 6: Two randomly selected retrieval results. The top row shows top-5 retrieved aerial images in descending order. The bottom row is the query sequence. The aerial images with blue border are the ground truth.

noticing that the second, third and fourth aerial images share most of their appearance with the top-1 image. In the bottom two rows of Fig. 6, although our model fails to predict the ground truth in the first place, we can see that visually the top-1 prediction is very similar to the ground truth.

### 5.2. Ablation Studies

To evaluate the effectiveness of our proposed models, we conducted ablation experiments. We study the effectiveness of the TFAM modules, SD, and number of heads in the TFAM module as reported in Table 3 and Table 4. We observe that the model which has 6 TFAMs with 8 heads for each multi-head self-attention layer and random dropout of a maximum of 6 images achieved the best results among all the configurations. Moreover, we evaluate performance of our model under different backbones in Table 5. To be no-

T	H	R@1	R@5	R@10	R@1%
0	0	0.91%	4.49%	7.98%	26.69%
2	2	1.45%	6.22%	10.02%	31.84%
4	2	1.40%	6.34%	10.31%	32.97%
4	4	1.51%	6.27%	10.51%	32.93%
6	4	1.59%	6.02%	9.88%	32.14%
6	8	<b>1.80%</b>	<b>6.45%</b>	<b>10.36%</b>	<b>34.38%</b>

Table 3: Ablation study on the number of heads and TFAMs. ‘T’ is short for number of TFAMs and ‘H’ is short for number of Heads.  $J$  is fixed to 6.

	R@1	R@5	R@10	R@1%
$J = 1$	1.40%	6.08%	9.45%	31.89%
$J = 3$	1.51%	6.64%	10.57%	34.34%
$J = 5$	1.63%	6.41%	10.49%	34.40%
$J = 6$	<b>1.80%</b>	<b>6.45%</b>	<b>10.36%</b>	<b>34.38%</b>

Table 4: Ablation study on the maximum number of masked frames  $J$ . The model is fixed to 6 TFAMs with 8 heads.

BackBone	R@1	R@5	R@10	R@1%
VGG16 [33]	1.80%	6.54%	10.36%	34.38%
ResNet18 [15]	1.58%	5.98%	10.14%	33.83%
ResNet34 [15]	1.71%	7.01%	11.67%	38.16%
ResNet50 [15]	2.07%	8.12%	13.16%	40.10%

Table 5: Comparison between different backbones of the proposed model.

ticed, our model with ResNet50 [15] can achieve 40% on R@1%. But to keep a fair comparison with baseline methods, we still use VGG16 [33] as the backbone.

### 5.3. Variant Sequence Lengths

In real-world scenarios, the ground-level sequences could have different numbers of images. Given that our model has been trained with the SD scheme, in this experiment, we vary the number of ground-level images in a sequence at the inference time by modifying the value of the SD mask  $A$ . We compare our model with and without SD. To simulate the worst possible real-world scenarios, we started by dropping the first 6 images and only leaving the last 1 image in the sequence as the last image has the smallest visible overlap area with the aerial image. We then test by dropping the first 4 images and 2 images respectively. The results in Fig. 7 demonstrated that our model with SD outperformed the model trained without it in most of the cases which proves that SD improves the model performance and feature coherency on variable length sequences.

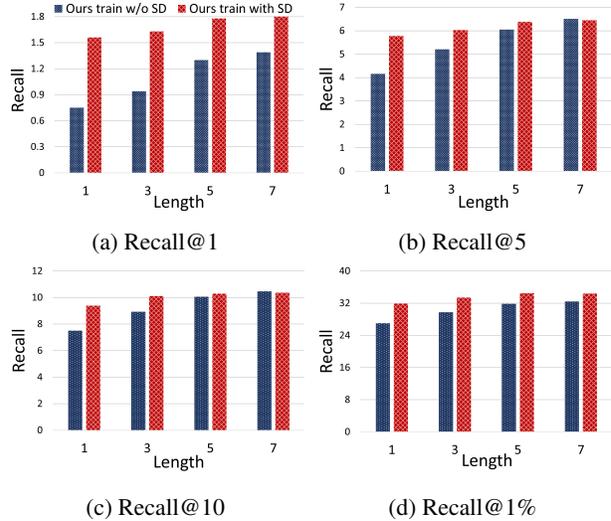


Figure 7: Comparison between variant sequence lengths simulated by SD in the testing phase for two models trained with (red) SD and without (black) SD. The results demonstrate that although trained using fixed-length sequences, the proposed SD enables our method to predict coherent feature representation compared with training without SD.

Noticeably, even ground sequence has only one single image during testing, our model trained with SD significantly outperforms SAFA [30] which was trained with the center image as query (Table 2) as observed from Fig 7.

## 6. Conclusion, Limitation, and Future Work

In this paper, we put forward the first cross-view geo-localization method that operates on sequences of limited FOV images. To aggregate the temporal features, we proposed a TFAM module that leveraged the multi-head self-attention mechanism to fuse information from a sequence of images. Although we used fixed-length sequences during the training phase, we simulated variant length sequences using our proposed sequential dropout method that regularizes our model to have a coherent feature representation. This also helps our model to tackle ground-level sequences with different lengths during the testing phase. We contributed to the vision community a novel large-scale cross-view *sequence* geo-localization dataset. Our extensive experiments demonstrated the effectiveness of different components of the proposed approach, robustness on variable-length input sequences, and state-of-the-art results against several competitive cross-view geo-localization methods.

One limitation of our proposed method is that the maximum length of the ground-level image sequence is constrained by the size of aerial images. Exploring methods that can geo-localize long sequences spanning multiple aerial images is one future research direction.

## References

- [1] Google maps static api. <https://developers.google.com/maps/documentation/maps-static/overview>.
- [2] Mapillary. <https://www.mapillary.com/app>.
- [3] Relja Arandjelović, Petr Gronat, Akihiko Torii, Tomas Pa-jdla, and Josef Sivic. Netvlad: Cnn architecture for weakly supervised place recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(6):1437–1451, 2018.
- [4] Joydeep Biswas and Manuela Veloso. Depth camera based localization and navigation for indoor mobile robots. In *RGB-D Workshop at RSS*, volume 2011, 2011.
- [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Matusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020.
- [6] Sudong Cai, Yulan Guo, Salman Khan, Jiwei Hu, and Gongjian Wen. Ground-to-aerial image geo-localization with a hard exemplar reweighting triplet loss. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- [7] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, pages 213–229, Cham, 2020. Springer International Publishing.
- [8] Francesco Castaldo, Amir Zamir, Roland Angst, Francesco Palmieri, and Silvio Savarese. Semantic cross-view matching. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV) Workshops*, December 2015.
- [9] Han-Pang Chiu, Varun Murali, Ryan Villamil, G. Drew Kessler, Supun Samarasekera, and Rakesh Kumar. Augmented reality driving using semantic geo-registration. In *2018 IEEE Conference on Virtual Reality and 3D User Interfaces (VR)*, pages 423–430, 2018.
- [10] Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. In *2005 IEEE computer society conference on computer vision and pattern recognition (CVPR’05)*, volume 1, pages 886–893. Ieee, 2005.
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009.
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- [14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [16] Sixing Hu, Mengdan Feng, Rang MH Nguyen, and Gim Hee Lee. Cvm-net: Cross-view matching network for image-based ground-to-aerial geo-localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7258–7267, 2018.
- [17] Dong-Ki Kim and Matthew R. Walter. Satellite image-based localization via learned embeddings. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2073–2080, 2017.
- [18] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.
- [19] Tsung-Yi Lin, Serge Belongie, and James Hays. Cross-view image geolocation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2013.
- [20] Tsung-Yi Lin, Yin Cui, Serge Belongie, and James Hays. Learning deep representations for ground-to-aerial geolocation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- [21] Liu Liu and Hongdong Li. Lending orientation to neural networks for cross-view geo-localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [22] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Ro{bert}a: A robustly optimized {bert} pretraining approach, 2020.
- [23] David Nistér, Oleg Naroditsky, and James Bergen. Visual odometry. In *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004.*, volume 1, pages I–I. Ieee, 2004.
- [24] Aude Oliva and Antonio Torralba. Modeling the shape of the scene: A holistic representation of the spatial envelope. *International Journal of Computer Vision*, 42(3):145–175, 2001.
- [25] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming

- Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems* 32, pages 8024–8035. Curran Associates, Inc., 2019.
- [26] Krishna Regmi and Mubarak Shah. Bridging the domain gap for ground-to-aerial image matching. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- [27] Krishna Regmi and Mubarak Shah. Video geo-localization employing geo-temporal feature learning and gps trajectory smoothing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 12126–12135, October 2021.
- [28] Eli Shechtman and Michal Irani. Matching local self-similarities across images and videos. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2007.
- [29] Akshay Shetty and Grace Xingxin Gao. Uav pose estimation using cross-view geolocalization with satellite imagery. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 1827–1833, 2019.
- [30] Yujiao Shi, Liu Liu, Xin Yu, and Hongdong Li. Spatial-aware feature aggregation for image based cross-view geo-localization. *Advances in Neural Information Processing Systems*, 32:10090–10100, 2019.
- [31] Yujiao Shi, Xin Yu, Dylan Campbell, and Hongdong 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 (CVPR)*, June 2020.
- [32] Yujiao Shi, Xin Yu, Liu Liu, Tong Zhang, and Hongdong Li. Optimal feature transport for cross-view image geo-localization. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(07):11990–11997, Apr. 2020.
- [33] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.
- [34] Yicong Tian, Chen Chen, and Mubarak Shah. Cross-view image matching for geo-localization in urban environments. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [35] Aysim Toker, Qunjie Zhou, Maxim Maximov, and Laura Leal-Taixe. Coming down to earth: Satellite-to-street view synthesis for geo-localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6488–6497, June 2021.
- [36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [37] Nam N. Vo and James Hays. Localizing and orienting street views using overhead imagery. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, pages 494–509, Cham, 2016. Springer International Publishing.
- [38] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and Huaxia Xia. End-to-end video instance segmentation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8741–8750, June 2021.
- [39] Daniel Wilson, Xiaohan Zhang, Waqas Sultani, and Safwan Wshah. Visual and object geo-localization: A comprehensive survey. *arXiv preprint arXiv:2112.15202*, 2021.
- [40] Scott Workman, Richard Souvenir, and Nathan Jacobs. Wide-area image geolocalization with aerial reference imagery. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [41] Menghua Zhai, Zachary Bessinger, Scott Workman, and Nathan Jacobs. Predicting ground-level scene layout from aerial imagery. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [42] Sijie Zhu, Taojiannan Yang, and Chen Chen. Revisiting street-to-aerial view image geo-localization and orientation estimation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 756–765, January 2021.
- [43] Sijie Zhu, Taojiannan Yang, and Chen Chen. Vigor: Cross-view image geo-localization beyond one-to-one retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3640–3649, June 2021.

## Supplementary Material

In this supplementary material, we are providing additional information for the following items:

- The availability of panoramas and limited Field-Of-View (FOV) images.
- Dataset coverage map
- More implementation details
- More details about the baseline methods discussed in the main paper.
- Comparison of the number of trainable parameters between the proposed model and baseline methods
- The availability of our proposed dataset and code for the public.
- More samples from our proposed dataset.
- More qualitative results predicted by our proposed model.

### A. Panorama vs Limited FOV images

As we discussed in our main paper, limited FOV images are more popular and common than panoramas. To highlight the difference, we presented the coverage areas of limited FOV images and panoramas from Mapillary [2] in Fig 8. Mapillary [2] is one of the largest crowdsourcing platforms for sharing geotagged photos. As of 2018, Mapillary [2] hosted 422 million images across the world. As observed from Fig 8, the coverage area of limited FOV images (Fig. 8b) on Mapillary is substantially greater than the coverage area of panoramas (Fig. 8a), especially in some developing areas such as Middle East, Africa and south America. We refer this to the complexity of capturing panoramic images which they need special and expensive cameras. To this end, using sequences of limited FOV images as the query is much more practical than using panoramas as the query in cross-view geo-localization.

### B. More implementation details

Our model was trained in an end-to-end manner using Adam [18] with weight decay of  $10^{-6}$  for 50 epochs on a single Nvidia V100 GPU. The learning rate is set initially to  $10^{-5}$  and decayed linearly to  $5 \times 10^{-7}$  after 30 epochs. We set the  $\gamma$  in Equation 5 of main paper to 10. We set the ground sequence length  $T = 7$  which is suitable for our dataset. We used the exhaustive mini-batch strategy [37] to construct the triplet pair with batch size set to 24.

## C. Baseline Methods

We employed two baseline methods for comparison, SAFA [30] and VIGOR [43]. For SAFA [30], we adopted their original code.<sup>3</sup> SAFA trained only on the center images of each sequence. For fair comparison, SAFA has been initialized with weights pretrained on CVUSA [40] dataset then trained on our dataset. We used same hyperparameters reported in SAFA’s original paper [30] and fine-tuned the model for 10 epochs. For VIGOR [43], we used their code<sup>4</sup> for training. Similar to SAFA, we trained their model from all images in the sequences by setting the center ground-level image to a ‘positive’ sample and the others are ‘semi-positive’ samples as defined in their original paper. We set the hyperparameters as reported in original VIGOR paper [43] and followed their exact procedures for training.

## D. Dataset Availability and Anonymity

Our proposed dataset is composed of two parts, ground-level image sequences and satellite imagery as explained in the main paper. Our ground-level images are public images collected by Vermont Agency of Transportation<sup>5</sup>. The private information of all ground-level images has been anonymized. These images will be shared publicly. Our satellite images came from Google Maps. Following Google Maps Platform Terms of Service<sup>6</sup>, we will make our dataset available for research purposes only. We will follow existing datasets, such as VIGOR [43], to distribute the collected dataset upon request.

## E. Dataset Coverage Map

To better visualize the diversity of the proposed dataset, we visualize the coverage area in Fig. 9. As indicated by the coverage map, our dataset includes both suburban and urban areas in Vermont, US which cover most scenarios on the roads.

## F. Comparison of parameters

In this section, we present the comparison of trainable parameters between the proposed model with different backbones and baseline methods in Table 6. Our model with VGG16 [33] is larger than the baselines. This is because the output dimension of VGG16 is 4096. As a result, we need wider TFAMs to handle this large latent vector. When we switch to ResNet [15] as backbone, the number of parameters is significantly less than VGG [33] as backbone. This is because the dimension of output of ResNet50

<sup>3</sup>[https://github.com/shiyujiao/cross\\_view\\_localization\\_SAF](https://github.com/shiyujiao/cross_view_localization_SAF)

<sup>4</sup><https://github.com/Jeff-Zilence/VIGOR>

<sup>5</sup><https://vtrans.vermont.gov/>

<sup>6</sup><https://cloud.google.com/maps-platform/terms>

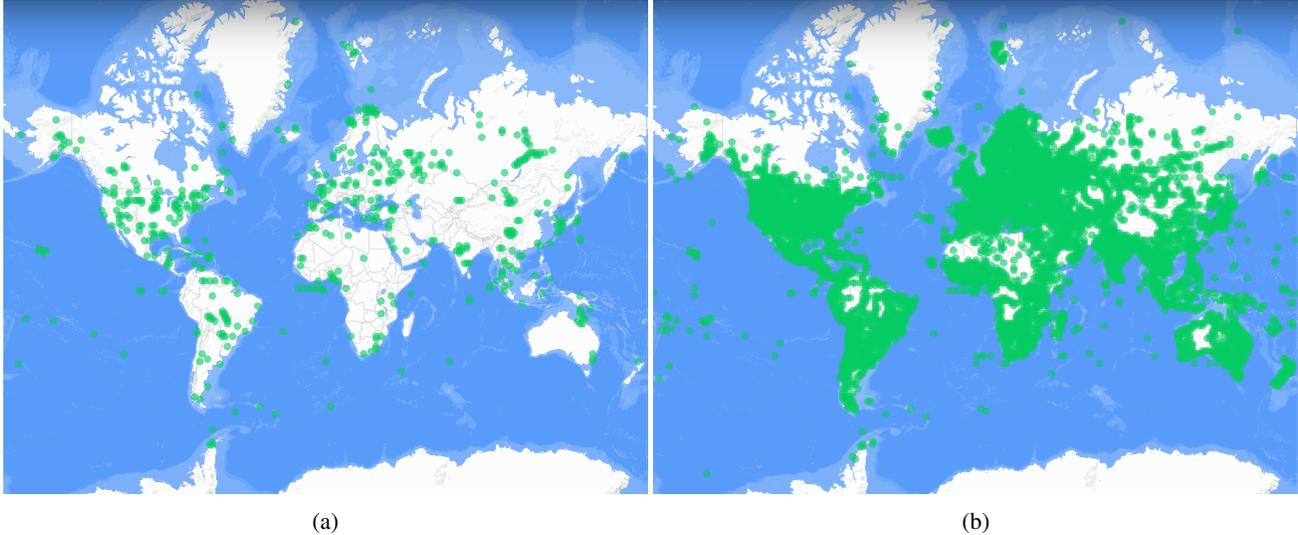


Figure 8: Comparison of coverage area (green lines) of user uploaded street view images between panoramic (a) and limited FOV images (b) on Mapillary [2].

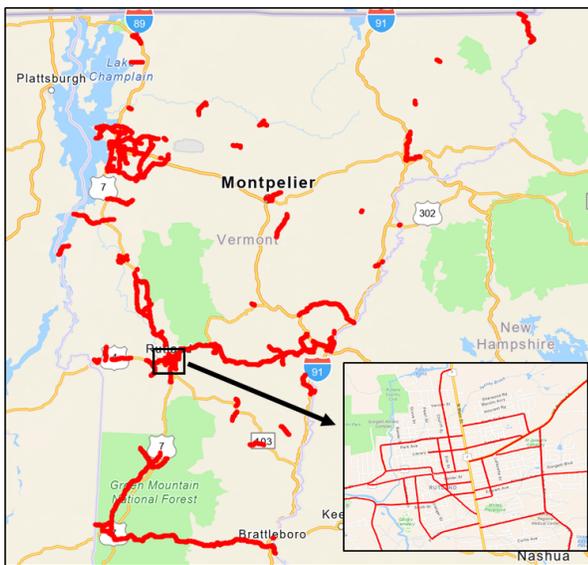


Figure 9: The coverage map of the proposed dataset. The coverage area is indicated by red lines.

is 2048. For ResNet34 and ResNet18, the dimension of output is only 512 which cause these two models are even smaller than baselines. However, despite of the backbones, the proposed model is constantly outperforms the baseline methods. For a fair comparison with baseline methods, we choose VGG16 as the backbone in the main script.

Method	Parameters	R@1	R@10	R@1%
VIGOR [43]	395M	0.54%	4.48%	18.55%
SAFA <sup>†</sup> [30]	319M	0.68%	5.06%	21.81%
Ours w/ VGG16 [33]	2.9G	1.80%	10.36%	34.38%
Ours w/ Res50 [15]	775M	2.07%	13.16%	40.10%
Ours w/ Res34 [15]	240M	1.71%	11.67%	38.16%
Ours w/ Res18 [15]	161M	1.58%	10.14%	33.83%

Table 6: Comparison between our proposed methods with different backbones and baseline methods. <sup>†</sup> indicates testing on single center ground image as query.

## G. More Dataset examples

In this section, we provided 6 randomly sampled satellite and ground sequence pairs from our proposed dataset as shown in Fig. 10. As shown in Fig. 10, our dataset covers diverse locations, urban, suburban, and rural areas which we discuss in detail in our main script.

## H. More Qualitative Results

In this section, we provided more retrieval examples. Fig. 11 shows correct top-1 examples predicted by our model and Fig. 12 shows top-5 retrieval examples. Each figure shows pairs of satellite and ground images ordered from top to bottom. For each pair, the bottom row is the query ground-level sequence and the upper row is the predicted top-5 satellite images ranked in descending order from left to right. The satellite images with blue boarder are the ground truth.

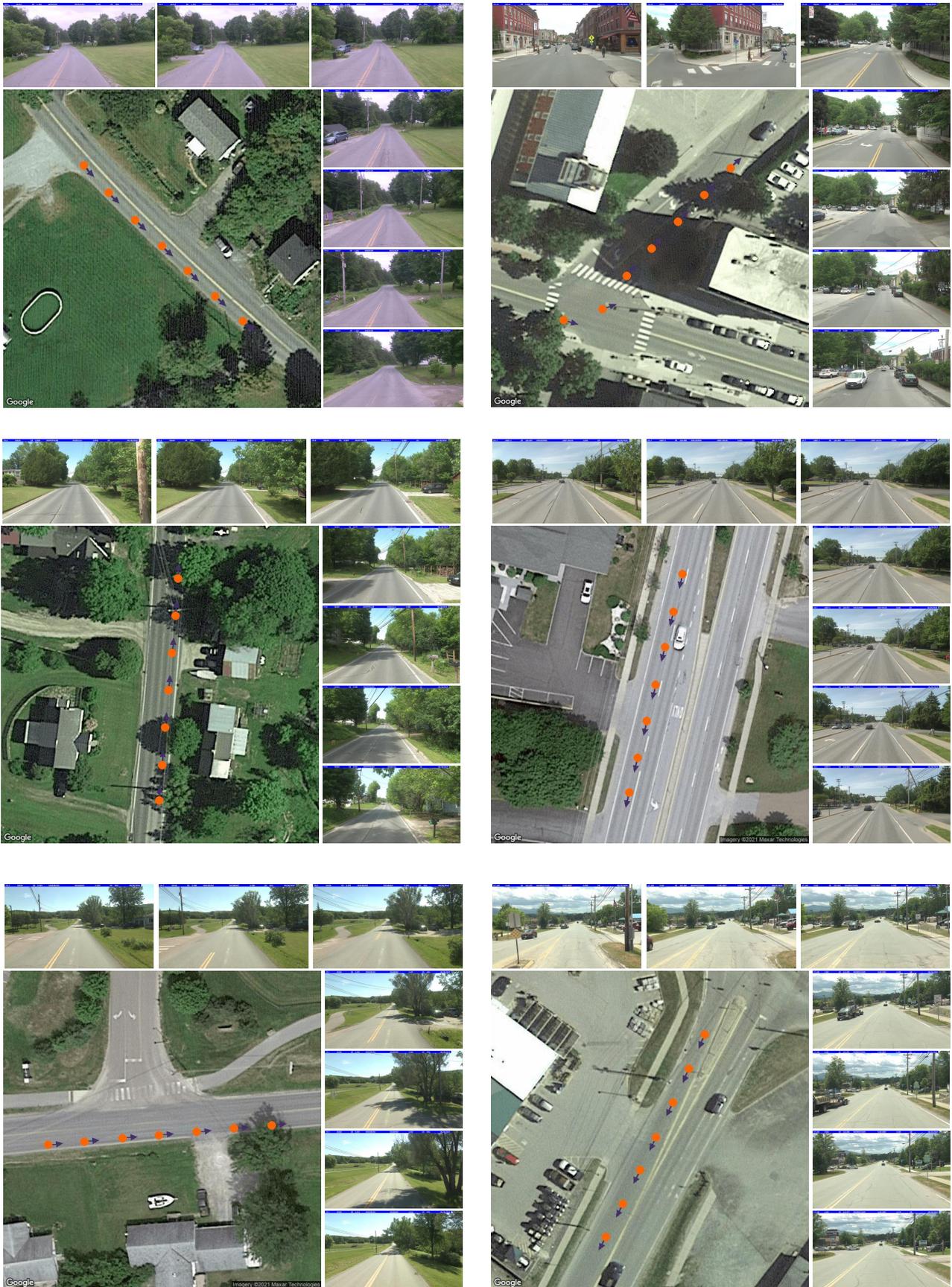


Figure 10: Six randomly sampled satellite and ground sequence pairs from our dataset.

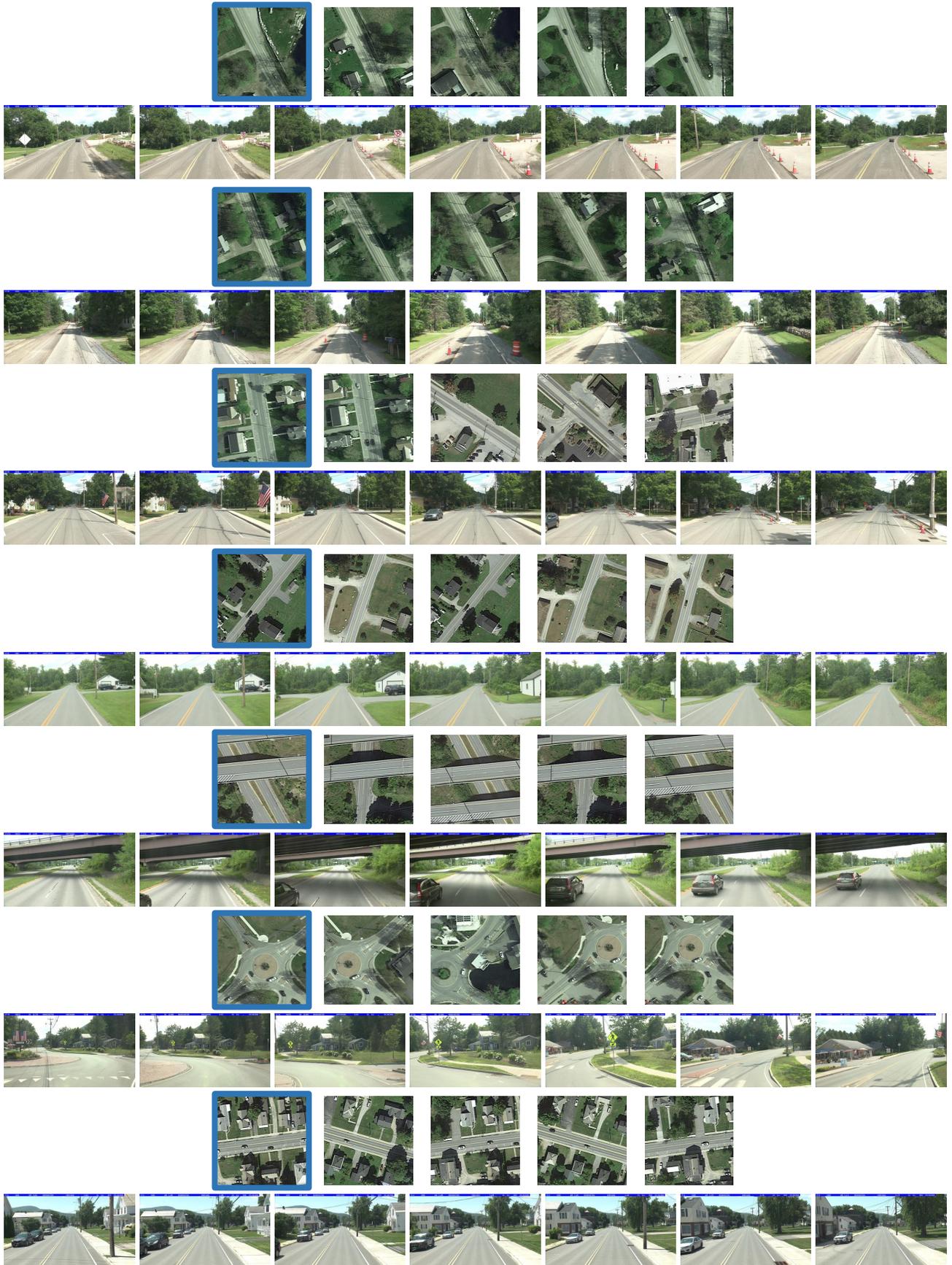


Figure 11: Samples been correctly predicted as top-1 by our model.



Figure 12: Samples been correctly predicted as top-5 by our model.