# Video2BEV: Transforming Drone Videos to BEVs for Video-based Geo-localization

Hao Ju
University of Macau
haojudalian@163.com

Zhedong Zheng*
University of Macau
zhedongzheng@um.edu.mo

## Abstract

*Existing approaches to drone visual geo-localization predominantly adopt the image-based setting, where a single drone-view snapshot is matched with images from other platforms. Such task formulation, however, underutilizes the inherent video output of the drone and is sensitive to occlusions and environmental constraints. To address these limitations, we formulate a new video-based drone geo-localization task and propose the Video2BEV paradigm. This paradigm transforms the video into a Bird's Eye View (BEV), simplifying the subsequent matching process. In particular, we employ Gaussian Splatting to reconstruct a 3D scene and obtain the BEV projection. Different from the existing transform methods, e.g., polar transform, our BEVs preserve more fine-grained details without significant distortion. To further improve model scalability toward diverse BEVs and satellite figures, our Video2BEV paradigm also incorporates a diffusion-based module for generating hard negative samples, which facilitates discriminative feature learning. To validate our approach, we introduce UniV, a new video-based geo-localization dataset that extends the image-based University-1652 dataset. UniV features flight paths at 30° and 45° elevation angles with increased frame rates of up to 10 frames per second (FPS). Extensive experiments on the UniV dataset show that our Video2BEV paradigm achieves competitive recall rates and outperforms conventional video-based methods. Compared to other methods, our proposed approach exhibits robustness at lower elevations with more occlusions.*

## 1. Introduction

Given the visual information captured by the drone, drone visual geo-localization aims to retrieve the image of the same location from another platform, *e.g.*, satellite, which is typically associated with the off-line GPS metadata [54]. This capability enables drones to locate themselves even in
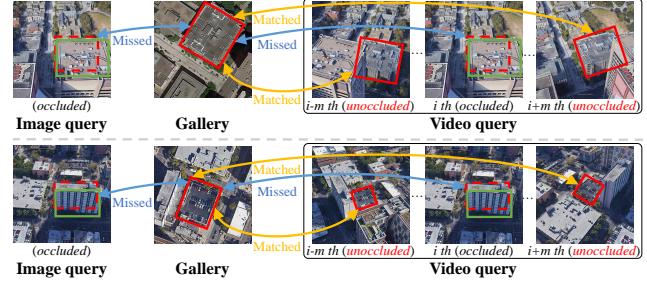
*Corresponding author.



Figure 1. **Typical failure cases for image-based drone geo-localization.** For image queries, the core areas in ground-truth are occluded by another building, largely compromising the spatial matching. In contrast, video queries usually contain unoccluded frames in a circling flight, and thus could reflect a more comprehensive view of the target location. Drone video inputs are relatively robust to occlusions.

GPS-unavailable areas, such as those between tall buildings or in rural regions. The common setting is based on the snapshot image from the drone as a query to search the matched location in the satellite-view candidate pool. Despite its potential, image-based drone geo-localization faces three primary challenges: (1) The single image from the drone view is often affected by occlusion and other environmental constraints. (2) The geometric transformation based on a single image usually introduces edge distortions and blurring. (3) Many locations share similar patterns, such as architectural styles, making it challenging to distinguish between them using a single image alone.

This paper attempts to address three challenges in drone visual geo-localization. First, existing approaches [5, 42, 43, 54] predominantly adopt the image-based setting, where a single drone-view snapshot is matched with images from other platforms. As shown in Fig. 1, such task formulation underutilizes the inherent video output of the drone and is sensitive to occlusions and environmental constraints. To mitigate this, we formulate a new video-based drone geo-localization task, which harnesses the drone video as input to capture the target location comprehensively. By controlling the drone to achieve unobstructed viewpoints, the video
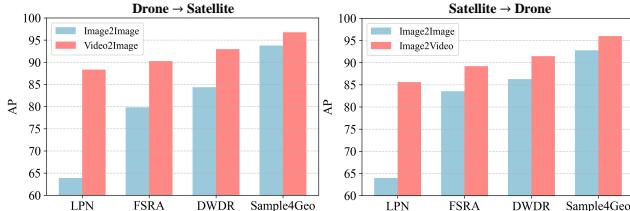
Figure 2. Performance comparisons of leveraging image data (Image2Image) or video data (Video2Image or Image2Video) with methods including LPN [42], FSRA [3], DWDR [43], and Sample4Geo [5]. We report the Average Precision (AP) metric. For a fair comparison, we keep the same number of data in the gallery. We could observe that all our re-implemented methods arrive at better performance when adopting video query or gallery.
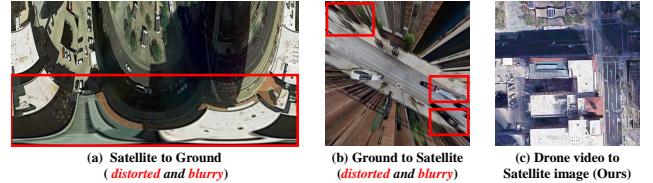


Figure 3. Prevailing image-based geometric transformation (a) Satellite to Ground transformation by the polar transformation [30], (b) Ground to Satellite transformation by the spherical transformation [45]. Our Drone video to Satellite transformation is shown in (c). Compared to image-based approaches, our method, fully leveraging the comprehensive view from the free-of-lunch drone videos, mitigates severe distortion and blurring.

input significantly reduces the impact of occasional occlusions present in the single frame, thereby improving the performance on all our re-implemented methods (see Fig. 2).

Second, geometric transformations based on a single image often introduce edge distortion and blurring (see Fig. 3). Existing approaches [30, 45] usually apply the transformation pre-processing to align the input images collected from different platforms. Such transformations coarsely align the images, thereby easing spatial comparison for deeply-learned models. For instance, Shi et al. [30] propose transforming satellite views to ground views via the polar transformation, which generally aligns the layout of the transformed images with fisheye ground-view photos. Similarly, Wang et al. [45] manipulate panoramic ground views to satellite views via the spherical transformation, simplifying the comparison with satellite-view data. However, both the pre-processing transformation methods suffer from spatial distortion and blurring, leading to incomplete image-level alignment. Since the input of our task is the video with different viewpoints along the flight, we propose a Video2BEV paradigm that transforms the video into a Bird's Eye View (BEV). Specifically, we employ Gaussian Splatting [12] to reconstruct a 3D scene and obtain the BEV projection. Unlike existing transformation methods, such as polar and spherical transformations, our transformation is not based on a single 2D image but a 3D scene derived from videos. Therefore, the core area of our BEVs preserve more fine-grained textures without significant distortion or blurring.

Third, many locations share similar patterns, such as architectural styles, making it challenging to distinguish between them using a single image alone. To address this, we resort to the drone video input in this work and conduct more discriminative feature learning. To enhance model scalability and discriminative capability, our Video2BEV paradigm incorporates a diffusion-based module for generating hard negative samples. This module generates appearance-similar BEVs, motivating the model to focus on fine-grained details to discriminate diffusion-generated BEV hard negatives and BEVs derived from videos.

In an attempt to overcome the aforementioned challenges, this paper (1) adopts the video setting, (2) transforms video to BEVs derived from reconstructed 3D scenes, (3) introduces a diffusion model to generate hard negatives. Given the lack of a video-based drone geo-localization benchmark, we introduce UniV, a new video-based geo-localization dataset that extends the image-based University-1652 dataset. UniV features flight paths at two different elevation angles with increased frame rates up to 10 frames per second (FPS). Extensive experiments on the UniV dataset show that our Video2BEV paradigm achieves competitive recall rates and outperforms conventional video-based methods. Compared to other methods, our proposed approach exhibits greater robustness at lower elevations with more occlusions. To summarize, the key contributions of our work are as follows:
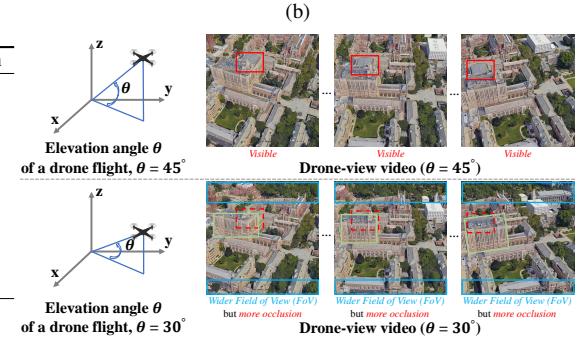
- We present a Video2BEV paradigm that transforms drone-view videos to Bird's Eye Views (BEVs), easing subsequent matching with satellite images. In particular, our Video2BEV introduces 3D Gaussian Splatting for geometric projection, while also including a hard negative generation module to learn from diverse BEVs and satellite images.

- Given the lack of a video-based drone geo-localization benchmark, we introduce a new video-based geolocalization dataset, UniV, with two elevation angles and 10 FPS. This dataset contains drone videos, satellite images, and ground images, and is closer to real-world scenarios containing typical cases, e.g., occlusions.

- Extensive experiments on the proposed UniV dataset show: (1) Image-based queries are highly sensitive to occlusions and environmental changes, such as occlusions, whereas video-based queries exhibit greater robustness (see Fig. 2). (2) The proposed Video2BEV, which incorporates geometric transformations, achieves 96.80 AP on Drone Video → Satellite, significantly outperforming other existing methods. We observe a similar result on the drone video with the lower elevation angle, which exposes more occlusions.

Table 1. (a) Dataset comparisons between UniV and other visual geo-localization datasets. G, S, and D denote ground-view, satellite-view and drone-view, respectively. We enable video modality and add another common elevation angle of drone flight. (b) Elevation angles $\theta$ illustration. Top panel shows $\theta = 45°$ and bottom panel displays $\theta = 30°$. With a lower elevation angle, the new flight captures the target building with *wider Field of View (FoV)* but more *occlusions*, thereby posing more challenges for drone visual geo-localization.

(a)

| Datasets | Platforms | #data/location | Modality | Elevation |
|---|---|---|---|---|
| CVUSA [47] | G, S | 1 image + 1 image | Image | N/A |
| Lin *et al*. [17] | G, S | 1 image + 1 image | Image | 45° |
| Vo *et al*. [41] | G, S | 1 image + 1 image | Image | N/A |
| Tian *et al*. [37] | G, S | 1 image + 1 image | Image | 45° |
| CVACT [18] | G, S | 1 image + 1 image | Image | N/A |
| Vigor [56] | G, S | 2 images + 1 image | Image | N/A |
| University-1652 [54] | G, S, D | (16 + 1 + 54) images | Image | 45° |
| GeoText-652 | G, S, D | (16 + 1 + 54) images + 180 texts | Image + Text | 45° |
| UniV | G, S, D | (16 +1) image + **2 videos** | **Video** | **30°**, 45° |

(b)

**Elevation angle $\theta$ of a drone flight, $\theta = 45°$** — Drone-view video ($\theta = 45°$) — *Visible* *Visible* *Visible*

**Elevation angle $\theta$ of a drone flight, $\theta = 30°$** — Drone-view video ($\theta = 30°$) — *Wider Field of View (FoV) but more occlusion* *Wider Field of View (FoV) but more occlusion* *Wider Field of View (FoV) but more occlusion*

# 2. Related Work

**Visual Geo-localization.** Visual geo-localization is to locate the position via the visual information, which is usually formulated as a sub-task of image retrieval [54]. The ancient people could check the surroundings to find their location via the paper map. Similarly, visual geo-localization searches the relevant position candidates recording in other platforms, *e.g*., satellite. The primary challenge of this task is due to the inherent difference among multiple platforms, *i.e*., appearance changes due to various viewpoints [30]. Therefore, previous methods focus on alignment to establish spatial correspondence between these platforms, which can be coarsely divided into two families, *i.e*., image-level alignment and feature-level alignment. For image-level alignment, Shi *et al*. [30] first propose to leverage polar transformation to warp satellite images to the ground view. Other methods [13, 45] transform panoramic ground images to the satellite view by proposed spherical transformation. Regmi *et al*. [26] utilize generative adversarial networks to synthesize aerial images for matching with ground images. However, the transformed images often suffer from severe distortion, blurring, and unrealistic content, which hinders fine-grained alignment between the two distinctive viewpoints. For implicit feature-level alignment, prior works [5, 21, 36, 50, 51, 54, 57] tend to rely on the power of neural networks or loss functions to align features between the query and gallery. Other methods focus on the explicit alignment of features from different views. Among these, some methods [18, 31, 32] encode orientation information and utilize it for feature-level alignment. Some methods propose to establish key-point [16, 35] or region alignment [3, 46, 48, 49] in the feature level. Previous works overlook the viewpoint variation within drone-view videos, treating them in an image-based setting. In this work, we aim to unlock the potential of videos captured by a drone and mitigate view disparity via the transformation without obvious distortion or blurring.

**Video Understanding.** The early works, *e.g*., Carreira *et al*. [2] and Shen *et al*. [29], leverage the two-stream 3D convolution network that combines video data with corresponding optical flow. In this way, motion information is adopted as the guidance for the fusion of different frames. Instead of adopting the 3D convolution, Feichtenhofer *et al*. [7] propose to deploy a two-branch 2D convolution network to fuse spatial semantic information and motion information together. Recently, attention mechanisms have shown their priority across many tasks [1, 4, 19, 40]. However, the computational cost of processing video data is substantial [8]. Many methods focus on designing efficient attention operations and training strategies. For efficient attention operations, Liu *et al*. [20] propose a 3D shifted window mechanism based on multi-head attention to fuse inherent complementary information of videos. Son *et al*. [34] propose CNN-based attention, modeling the ordering reliance in video frames and capturing the long-term dependency. For effective training strategies, Tong *et al*. [38] introduce mask image modeling [9] into video data in a self-supervised manner, leveraging complementary information across different frames. To keep task-specific and shared knowledge across different tasks, Peirone *et al*. [24] treat video understanding as a unified task and propose to train task-specific heads and a cross-task backbone simultaneously. Different from the daily videos, drone videos for geo-localization typically contain multi-view information for the target location [22, 44]. Therefore, rather than adopting the off-the-shelf video backbone, we propose a Video2BEV transformation to leverage the 3D geometric correspondences and enable a straightforward spatial alignment for matching.

# 3. The UniV Dataset

Given the lack of a video-based drone geo-localization benchmark, we collect a new dataset dubbed *UniV* involving the video modality. We follow the building in-

formation and the protocol of the existing University-1652 dataset [54]. The UniV dataset encompasses 1,652 buildings in 72 universities from three platforms, *i.e.*, ground, satellite, and drone cameras. In particular, the UniV dataset contains 16 ground-view images and 1 satellite-view image for each building and the training set of UniV dataset contains 701 buildings, while the test set in the UniV dataset includes other 951 buildings. There are no overlapping buildings between the training and test sets. The proposed UniV dataset is different from the image-based University-1652 and other datasets in two primary aspects, *i.e.*, modality and elevation-angle expansions (see Tab. 1a).

**Modality Expansion.** Existing datasets [18, 41, 47, 56] collect data from two platforms, *e.g.*, satellite and ground cameras. Although some datasets [17, 37, 54] include drone or aerial views, they still collect data in the image format. We adopt similar operations as the University-1652 dataset but collect drone-view data in video format. Specifically, we leverage the 3D engine of google earth [23] to simulate the movement of a drone-view camera in the real world. To collect video data containing both scale and viewing-point variations, we leverage the dynamic viewpoints within the 3D engine and set the moving viewpoints along a spiral curve for moving around the target building three circles, closely approximating real-world drone flights. Video data are collected in 30 frames per second. Considering the video redundancy, in practice, we subsample videos along the temporal dimension, resulting in frame rates of 2, 5, and 10 for further processing.

**Elevation-angle Expansion.** Conventional datasets [17, 37, 54] collect drone or aerial data in a fixed elevation angle, *i.e.*, $45°$, which does not fully simulate the real-world use cases. Therefore, we add one new synthetic flying path at another common setting, *i.e.*, a lower elevation angle $30°$. The new flying path poses two new challenges (see Tab. 1b). First, drones flying at a $30°$ elevation angle capture scenes that include the target building and more surrounding areas, providing a wider Field of View (FoV). At the same time, it introduces disruptions for the center target building during matching. Second, there are more occluded cases, which lay over the core areas of the target building. It poses challenging to mine the discriminative frames in the video, while it is easier in the same location captured at a $45°$ elevation. Therefore, the proposed dataset could further evaluate the robustness of the method against more disruptions and occlusions, which is closer to real-world drone usage.

**Discussion. The contribution to the community.** The key difference between the UniV dataset and existing datasets [17, 37, 54] lies in the modality expansion of videos (see Tab. 1a), facilitating the development of robust drone visual geo-localization. A single image provides limited information about the corresponding building. When core areas of the buildings are occluded, single-image queries can

not produce reliable matching results (see Fig. 1). In such cases, the video contains both occluded and unoccluded frames. One frame may contain core-area information to complement another frame and together they can provide robust and complete information required for drone visual geo-location. In this way, all re-implemented methods perform better when adopting video data (see Fig. 2). Moreover, the UniV dataset also introduces a new real-world challenge for drone visual geo-localization. The new elevation angle of $30°$ is typical in real-world flights*. The $30°$ elevation angle faces more occlusion cases (see Tab. 1b), simulating outputs of real-world drone flights.

## 4. Method

### 4.1. Video2BEV Transformation

During the flight around the target building, the viewpoints of the camera vary, resulting in captured drone-view videos that contain rich multi-view information about both the target building and surrounding areas. Rather than taking viewpoint variation as an obstacle for aligning satellite and drone views, we explicitly leverage the multi-view information and transform the drone-view video into Bird's Eye Views (BEVs). In doing so, we ease the learning process for the model. Instead of learning geometry correspondence and feature correspondences simultaneously, the model only needs to learn the feature mapping relationship between two views, thus significantly facilitating network convergence. Illustrated in the left of Fig. 4, given the drone-view video containing multi-view images, we estimate corresponding camera poses by structure from motion [39] and reconstruct the scene containing the target buildings utilizing 3D Gaussian Splatting (3DGS) [12]. After reconstructing the scene, we adopt the normalized input pose and the unit vector in the world coordinate to calculate the BEV camera pose and render BEVs. For the training set, we incorporate rotation angles and varying heights into the BEV camera pose, generating a sequence of rotating and scaled-down BEVs. In terms of the test set, we only increase the height of the BEV camera poses and render a sequence of scaled-down BEV images. The number of BEVs corresponds to the number of images in the drone video. As shown in Fig. 4, outputs of Video2BEV transformation do not suffer from severe distortion, thereby aiding in the establishment of fine-grained spatial correspondence with the satellite view.

### 4.2. Hard Negative Sample Synthesis

Negative samples play a significant role in discriminative metric learning. Current negative sample mining strategies cannot ensure the quantity and quality of negative samples,

---

*The United States and the United Kingdom allow drone flights up to 400 feet; China restricts drones up to 120 meters.
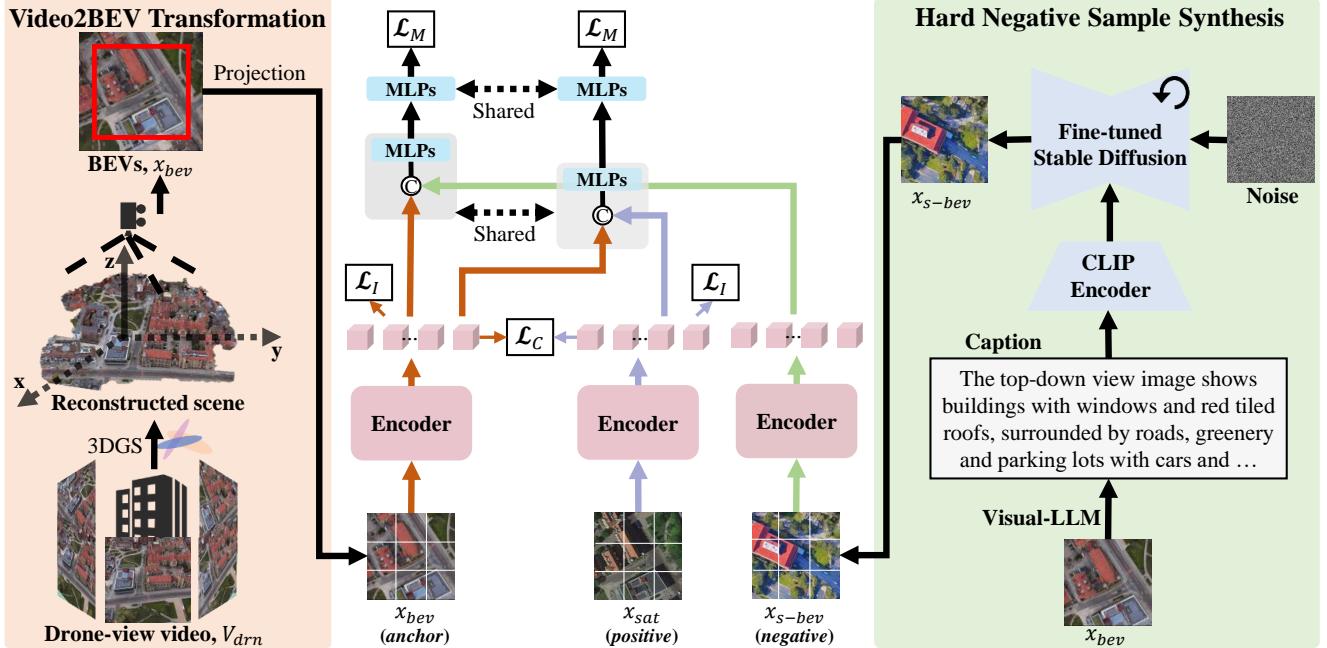
Figure 4. The overview of the Video2BEV paradigm. **Video2BEV Transformation (*left*).** Given drone-view video $V_{drn}$ containing multi-view frames, we adopt 3D Gaussian Splatting (3DGS) to reconstruct the scene at first. Then we render the scene from a Bird-Eye-View to get the projection (BEVs). Considering the region of the core area, we further crop BEVs for training. We can observe that BEVs exhibit resemblances to the corresponding satellite-view images. **Hard Negative Sample Synthesis (*right*).** Given captions generated by an off-the-shelf visual-LLM [11], we fine-tune a stable-diffusion model [28] with LoRA [10], and conduct inference to synthesize samples which serve as negative samples for subsequent usage. **Model Architecture (*middle*).** Given outputs of the proposed Video2BEV transformation, we extract embeddings by a shared encoder for satellite images $x_{sat}$ and BEVs $x_{bev}$, supervised by the contrastive loss $\mathcal{L}_C$ and the instance loss $\mathcal{L}_I$. Then we extract embeddings from synthetic BEVs $x_{s-bev}$ and adopt Multilayer Perceptrons (MLPs) to fuse both positive and negative samples, supervised by the matching loss $\mathcal{L}_M$. Similar operations are also applied for satellite-view images which are omitted.
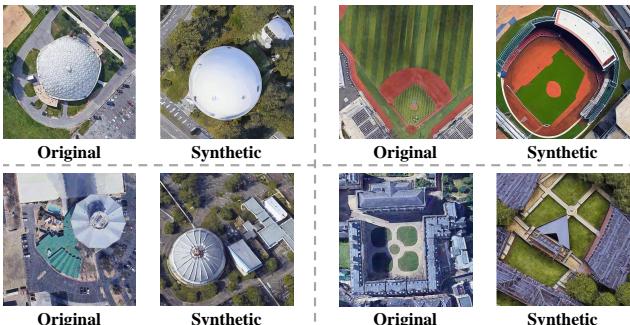


Figure 5. Visualizations of original images and synthetic hard negatives. Synthetic negatives exhibit similar colors and structures to original images, which assures the quality of negatives.

as the number of challenging samples is limited, and selected negative samples can not always share similar architectural styles and other details with the original samples. In order to bypass these drawbacks, we propose to fine-tune a diffusion model to generate diverse samples as negative samples, shown in the right part of Fig. 4. After transforming the drone-view video to BEVs $x_{bev}$ via the Video2BEV transformation, we utilize a visual-LLM [11] to generate captions in a multi-round process for both BEV and satellite-view images, during which a human annotator interacts with visual-LLM to restrict the token length of the captions. After obtaining captions for the BEV and satellite images, we fine-tune a stable diffusion network [28] with LoRA [10] to generate diverse synthetic images, during which we freeze the CLIP text encoder [25]. The outputs of this model are challenging negative samples for the subsequent step. Considering the transformed BEV and satellite images share the same viewing direction and content, we conduct inference on the same stable diffusion network to generate negative samples for BEV and satellite-view images by adopting corresponding captions. We provide visualizations of original and synthetic images in Fig. 5. Similar to the original samples, synthetic negative samples exhibit comparable appearances, including color patterns and structural features of buildings.

## 4.3. Model Optimization

In this section, we do not overstate the novelty of the architecture. Instead, we adopt a general architecture from vision-language models [14, 15], enhanced by the synthetic negative samples. The model architecture is shown in the

middle of Fig. 4 and is optimized in a two-stage manner. In the first stage, we transform the drone-view video to BEVs by the proposed Video2BEV transformation (see the left part of Fig. 4). Then, we adopt a shared encoder to extract embeddings from paired BEV and satellite-view images. The encoder is ViT-S [6] excluding the classifier. The supervisions of this stage are the instance loss $\mathcal{L}_I$ with the square-ring partition [42] and the contrastive loss $\mathcal{L}_C$ [14]. For the instance loss, we first partition the embeddings from different views into different parts adopting the square-ring partition [42] strategy. Next, we apply multiple classifier modules to each part of the embeddings (similar to LPN [42]), yielding the location probability. Then we accumulate instance losses from multiple parts, and the instance loss $\mathcal{L}_I$ [55] is formulated as the location classification:

$$\mathcal{L}_I = -log(\hat{p}_{sat}) - log(\hat{p}_{bev}), \qquad (1)$$

where $\hat{p}_{sat}$ and $\hat{p}_{bev}$ are the predicted probability that belongs to the ground-truth label from two views respectively. For the contrastive loss, given a pair of satellite-view and BEV images, the satellite-to-BEV similarity is defined as:

$$S_{sat2bev} = \frac{exp(s(f_{sat}, f_{bev})/\tau)}{\Sigma_{j=1}^{N} exp(s(f_{sat}, f_{bev}^{j})/\tau)}, \qquad (2)$$

where $f_{sat}$ and $f_{bev}$ are the embeddings of the same location from two platforms, and $f_{bev}^{j}$ denotes the sample within the mini-batch. $\tau$ is a learnable temperature parameter. $s(\cdot, \cdot)$ denotes the cosine similarity. Similarly, the BEV-to-satellite similarity is $S_{bev2sat}$ and the contrastive loss $\mathcal{L}_C$ is:

$$\mathcal{L}_C = -\frac{1}{2}(log(S_{sat2bev}) + log(S_{bev2sat})). \qquad (3)$$

For the second stage, we employ two-layer Multilayer Perceptrons (MLPs) alongside the square-ring partition [42] to fuse two embeddings derived from different views of the input including BEV, satellite, and synthetic BEV, and synthetic satellite view images. Then we concatenate all parts of embeddings from two views and project the fused embeddings into the two-dimensional space by another MLPs. The matching loss of two-view inputs $\mathcal{L}_M$ is:

$$\mathcal{L}_M = -(p_m log(\hat{p}_m) + (1 - p_m)log(1 - \hat{p}_m)), \qquad (4)$$

where $\hat{p}_m$ is the estimated matching probability and $p_m$ is a ground-truth binary label. If the two input data do not contain synthetic data and are from the same category, then $p_m = 1$; otherwise, $p_m = 0$. Specifically, for the BEVs, we calculate the matching loss two times. For the first calculation, we rank the similarity $S$ and select *three* negative samples from the satellite-view images, ensuring that these samples do not belong to the same category. For the second calculation, we similarly select *another three* negative

samples from the synthetic BEV images, which includes samples from the same category that are actually challenging negative samples. We apply similar operations for the satellite-view input as well. Finally, we accumulate and average matching losses across different combinations of inputs. In summary, the loss functions in our method include the instance loss $\mathcal{L}_I$, the contrastive loss $\mathcal{L}_C$, and the matching loss $\mathcal{L}_M$. Specifically, we train the first stage of our method (including the encoder, classifier modules, and the temperature parameter $\tau$) with the instance loss $\mathcal{L}_I$ and the contrastive loss $\mathcal{L}_C$ at first. Subsequently, we freeze the fine-tuned first-stage weights and train the second stage (including two types of MLPs) from scratch, supervised by the matching loss $\mathcal{L}_M$.

**Discussion. What are the advantages of the synthetic negative samples?** First, inspired similar success in other fine-grained tasks [33, 52, 53], we encourage the model "see" more samples to prevent over-fitting as well as facilitate discriminative feature learning. With the assistance of the diffusion model, we can synthesize negative samples of a diverse range of *categories* and *quantities*. In terms of *categories*, the synthetic samples exhibit similar yet different architectural styles as shown in Fig. 5, including variations in color and pattern, as well as structural details such as the shape and material. Regarding *quantities*, the negative sample pool is no longer constrained to a fixed size. Utilizing the diffusion model, we can theoretically generate an infinite number of images as negative samples, expanding the negative sample pool significantly. Second, there is a clear relationship between our synthetic negative samples and the anchor samples. This is because we use identical captions from the original samples to synthesize the negative samples. This relationship ensures that our negative samples are appropriately challenging.

## 5. Experiment

**Implementation Details.** For the synthetic negative samples, the captions for BEV and satellite images differ and we synthesize 32 negative samples for both BEV and satellite images of each category. We train the first stage of the proposed model with the AdamW optimizer, with a batch size of 140, for 140 epochs, and a learning rate of $2e^{-5}$ and $2e^{-4}$ for the encoder and other modules in the first stage respectively. Then we freeze parameters in the first stage and train the second stage from scratch with a similar training configuration. During the test stage, we utilize the similarity scores from the first stage of the proposed model to select the top 32 samples from the gallery, and then apply the second stage of the proposed model to re-rank these top 32 samples. The whole framework is implemented with Pytorch.

**Evaluation Metrics.** Satellite-view data is in image format, while drone-view data is collected in video format.

Table 2. Comparisons on the UniV dataset for geo-localization between Drone (D) and Satellite (S). All methods are compared in the video setting. R@1 is recall at top1. AP (%) is average precision (high is good). $\theta$ denotes the elevation angle of the drone flight. The proposed method yields the best results.

| Method | $\theta = 45°$ | | | | $\theta = 30°$ | | | |
| | D→S | | S→D | | D→S | | S→D | |
| | R@1 | AP | R@1 | AP | R@1 | AP | R@1 | AP |
|---|---|---|---|---|---|---|---|---|
| LPN [42] | 86.31 | 88.34 | 83.31 | 85.60 | 68.62 | 72.50 | 67.76 | 71.30 |
| FSRA [3] | 88.59 | 90.25 | 87.30 | 89.17 | 81.60 | 84.17 | 77.89 | 81.00 |
| DWDR [43] | 91.73 | 92.96 | 89.87 | 91.45 | 88.02 | 89.81 | 85.59 | 87.85 |
| Sample4Geo [5] | 96.29 | 96.75 | 95.29 | 95.99 | 83.02 | 86.00 | 80.45 | 82.68 |
| Ours | **96.29** | **96.80** | **96.01** | **96.57** | **91.73** | **93.01** | **92.58** | **93.65** |



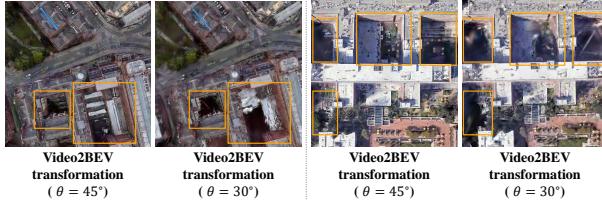| Video2BEV transformation ($\theta = 45°$) | Video2BEV transformation ($\theta = 30°$) | Video2BEV transformation ($\theta = 45°$) | Video2BEV transformation ($\theta = 30°$) |

Figure 6. The transformed BEV comparison of videos with different evaluation $\theta$. We highlight the challenging regions.
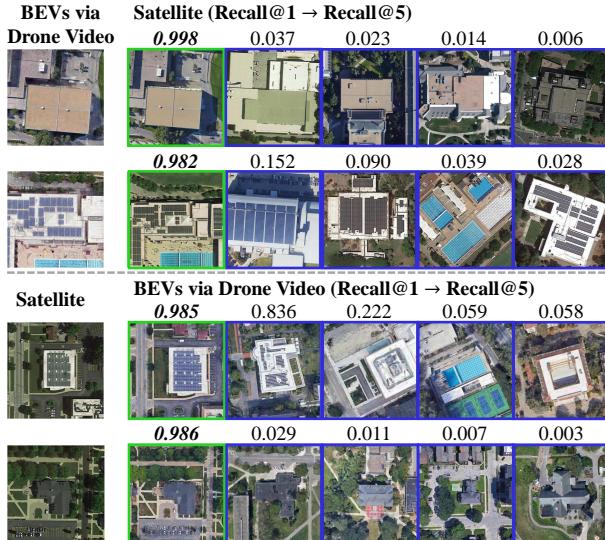


Figure 7. Qualitative results for Drone → Satellite and Satellite → Drone. We replace the visualization of the drone-view videos with the transformed output (BEVs) as the query or gallery. Given queries (left) from different platforms, matched galleries are in green box, and mismatched galleries are in blue box. The scores on the top are similarity scores estimated by the proposed method.

We can treat drone view data as images or video. In this paper, we adopt the video setting for the evaluation of competitive methods and our method. Specifically, we treat a drone video as an individual query or gallery by averaging the similarity scores of the images within the video in a late fusion manner. In our method, there is a similar averaging operation on similarity scores of the BEVs which is also in video format. We employ 2-fps videos in UniV as the training and test sets and release videos at 5 and 10 fps for further research usage.

## 5.1. Comparisons with Competitive Methods

**Quantitative Results.** As shown in Tab. 2, we compare the proposed method with other competitive methods on the UniV dataset. The performance of our method has surpassed that of other competitive methods [3, 5, 42, 43]. On the 45° subset, our method achieves gains of 0.30% Recall and 0.58% AP for satellite → drone compared to the second-best method. On the 30° subset, all methods experience a performance drop. As shown in Fig. 6, we highlight some imperfect reconstructed regions by the Video2BEV transformation. The lower elevation of the drone flights raises more occlusions (see Tab. 1b), which also compromises our Video2BEV transformation. Compared to the second best method, our method is still robust, receiving improvements of 3.2% AP for drone → satellite and 5.8% AP for satellite → drone respectively (see Tab. 2). All methods are compared in the video setting, which means we temporally average the outputs of frames in a video from the drone view. For methods with officially released weights (Sample4Geo, DWDR), we test these methods on the 45° test set directly and subsequently retrain and evaluate these methods on the 30° subsets. For methods without official weights (LPN, FSRA), we retrain these methods on both the 45° and 30° subsets to ensure a fair comparison.

**Qualitative Results.** We show qualitative results of the drone geo-localization in Fig. 7. In our method, drone-view videos are transformed to BEVs by the proposed Video2BEV transformation and we choose the representative sample from the BEV sequence for visualization. For drone → satellite, we observe that the proposed method effectively retrieves reasonable buildings with similar structural features, such as cross-shaped roofs and roofs equipped with solar panels. For satellite → drone, we find a similar result. Our method successfully retrieves true-matched results at the top of the candidate list among images with similar contents. *We add more qualitative visualizations including failure case analysis in the supplementary material.*

## 5.2. Ablation Studies and Further Discussion

**Effect of Primary Components.** We conduct ablation studies on the UniV dataset (45° subset). We employ the first stage of our method as the baseline (**Baseline**), which consists of a shared backbone supervised with the instance loss and contrastive loss. The input data for the baseline are drone-view videos and satellite-view images. Then, we transform drone-view videos to BEVs via the proposed Video2BEV transformation and adopt BEVs as input for the baseline, denoted as **BEV**. Next, we introduce the second stage of our method to the baseline, which is supervised by the matching loss, denoting **Two Stage**. The negative samples for this architecture are from in-batch samples [14]. Finally, we incorporate the synthetic negative samples in

Table 3. Ablation studies on: (a) Video2BEV transformation, the second stage of our method, and synthetic negative samples. (b) Different training strategies. **Train Together**: we fine-tune the first stage based on the weights pre-trained on ImageNet [27], and train the second stage from scratch at the same time. **Fine-tune**: we load fine-tuned first-stage weights on UniV, and then train both the first stage and the second stage. **Freeze**: we load fine-tuned first-stage weights on UniV, then fix the first-stage weights and only train the second stage from scratch. Notably, the **Freeze** strategy yields the best results. (c) Re-ranking different top-k samples in the second stage of our method. Considering the balance between performance and testing time, we choose to re-rank top-32 samples. D and S denote Drone and Satellite respectively. R@1 is recall at top1. AP (%) is average precision (high is good).

(a)

| Method | Vidoe2BEV transformation | Second stage | Synthetic negatives | D→S R@1 | D→S AP | S→D R@1 | S→D AP |
|---|---|---|---|---|---|---|---|
| Baseline | ✗ | ✗ | ✗ | 89.87 | 91.28 | 90.01 | 91.36 |
| BEVs | ✓ | ✗ | ✗ | 95.01 | 95.64 | 93.44 | 94.44 |
| Two Stage | ✓ | ✓ | ✗ | 95.86 | 96.48 | 95.01 | 95.78 |
| Ours | ✓ | ✓ | ✓ | **96.29** | **96.80** | **96.01** | **96.57** |

(b)

| Strategy | Load fine-tuned first-stage weights | Train first stage | Train second stage | D→S R@1 | D→S AP | S→D R@1 | S→D AP |
|---|---|---|---|---|---|---|---|
| Train Together | ✗ | ✓ | ✓ | 74.75 | 79.29 | 82.17 | 85.39 |
| Fine-tune | ✓ | ✓ | ✓ | 96.29 | **96.83** | 95.29 | 95.99 |
| Freeze | ✓ | ✗ | ✓ | **96.29** | 96.80 | **96.01** | **96.57** |

(c)

| Top-K | D→S R@1 | D→S AP | S→D R@1 | S→D AP |
|---|---|---|---|---|
| 8 | 96.01 | 96.52 | 95.58 | 96.10 |
| 16 | 96.01 | 96.51 | 95.72 | 96.25 |
| 32 | 96.29 | 96.80 | **96.01** | 96.57 |
| 64 | 96.29 | 96.81 | **96.01** | **96.60** |
| 128 | **96.43** | 96.98 | **96.01** | **96.60** |
| 256 | **96.43** | **96.99** | **96.01** | **96.60** |
| 512 | **96.43** | **96.99** | **96.01** | **96.60** |

Sec. 4.2 to train the second stage of our method and form the final version of our method, referred to as **Ours**. As shown in Tab. 3a, BEVs receive the largest performance improvement. We attribute this improvement to the reduction of the appearance gap between the drone-view images and the satellite-view images through the proposed Video2BEV transformation. Additionally, synthetic negative samples contribute to a substantial performance boost due to the enhanced quality of the negative samples for the second stage of Ours. The two-stage method (Two Stage) also receives improved performance, indicating that many false negative predictions are ranked within the range of the top 32. A fine-grained re-ranking can effectively rectify the matching results from the first stage of our method.

**Effect of Training Strategies.** We explore three different strategies for training the proposed model. For the **Train Together** strategy, we load the matched weights pre-trained on the ImageNet dataset [27]. Then we fine-tune the first stage of the proposed model and train the second stage of the proposed model from scratch simultaneously. The **Fine-tune** strategy entails loading fine-tuned weights of the first stage on the UniV dataset. After this, we fine-tune the first stage with a smaller learning rate while simultaneously training the second stage from scratch. The **Freeze** strategy consists of loading fine-tuned weights of the first stage on UniV, then fixing all weights of the first stage, while training the second stage from scratch at the same time. The results of three training strategies are in Tab. 3b. The **Train Together** strategy yields the worst results. We attribute this to the difficulty of training both stages simultaneously, as the first stage of the proposed model is designed for coarse-grained retrieval, while the second stage of the proposed model focuses on fine-grained retrieval, which relies on the output of the first stage. When both stages are trained together, the first stage fails to retrieve reliable candidates for the second stage, affecting the overall training process.

The **Fine-tune** strategy achieves a significant performance boost, as the first stage is able to produce reliable embeddings for the second stage. Finally, we freeze the first stage after loading its corresponding weight. The **Freeze** strategy receives the best result, and we adopt this training strategy. **Effect of Re-ranking Top-K Samples.** During the test stage, we select top-k samples from the gallery leveraging the similarity score from the first stage of our method and subsequently re-rank these samples by the second stage. We conduct hyper-parameter experiments with varying values of top-k. As shown in Tab. 3c, we select k ∈ {8, 16, 32, 64, 128, 256, 512}. Re-ranking the top-512 and top-256 samples yields the best performance and re-ranking the top-256, top-128, top-64, and top-32 samples results in a slight performance drop respectively. Re-ranking the top-16 and top-8 samples leads to a further decline in performance. Considering the balance between the performance of our method and the testing time, we choose to re-rank the top 32 samples as default.

# 6. Conclusion

In this work, we propose to leverage videos to mitigate the impact of environmental constraints in drone visual geo-localization. We propose a new Video2BEV paradigm that transforms drone-view videos into Bird's Eye View (BEV) images by 3D gaussian splatting. This transformation effectively reduces the viewpoint gap between the drone view and the satellite view. Our Video2BEV paradigm also includes a diffusion-based module to generate negative samples, enhancing the scalability of the model. To validate the proposed method, we introduce the UniV dataset, a new video-based drone geo-localization dataset. The dataset includes flight paths of the drone at $30°$ and $45°$ elevation angles and corresponding videos recorded at up to 10 frames per second. Our Video2BEV paradigm outperforms other approaches in terms of Recall@1 and AP.

# References

[1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *ICCV*, pages 6836–6846, 2021. 3

[2] João Carreira and Andrew Zisserman. Quo vadis, action recognition? A new model and the kinetics dataset. In *CVPR*, pages 4724–4733, 2017. 3

[3] Ming Dai, Jianhong Hu, Jiedong Zhuang, and Enhui Zheng. A transformer-based feature segmentation and region alignment method for uav-view geo-localization. *IEEE TCSVT*, 32(7):4376–4389, 2021. 2, 3, 7

[4] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers. *arXiv:2309.16588*, 2023. 3

[5] Fabian Deuser, Konrad Habel, and Norbert Oswald. Sample4geo: Hard negative sampling for cross-view geo-localisation. In *ICCV*, pages 16847–16856, 2023. 1, 2, 3, 7

[6] 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. *ICLR*, 2021. 6

[7] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *ICCV*, pages 6202–6211, 2019. 3

[8] Kai Han, Yunhe Wang, Hanting Chen, Xinghao Chen, Jianyuan Guo, Zhenhua Liu, Yehui Tang, An Xiao, Chunjing Xu, Yixing Xu, et al. A survey on vision transformer. *IEEE TPAMI*, 45(1):87–110, 2022. 3

[9] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, pages 16000–16009, 2022. 3

[10] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *ICLR*, 2022. 5

[11] Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, et al. Minicpm: Unveiling the potential of small language models with scalable training strategies. *arXiv:2404.06395*, 2024. 5

[12] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM TOG*, 42(4):139–1, 2023. 2, 4

[13] Guopeng Li, Ming Qian, and Gui-Song Xia. Unleashing unlabeled data: A paradigm for cross-view geo-localization. In *CVPR*, pages 16719–16729, 2024. 3

[14] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *NeurIPS*, 34:9694–9705, 2021. 5, 6, 7

[15] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, pages 12888–12900, 2022. 5

[16] Jinliang Lin, Zhedong Zheng, Zhun Zhong, Zhiming Luo, Shaozi Li, Yi Yang, and Nicu Sebe. Joint representation learning and keypoint detection for cross-view geo-localization. *IEEE TIP*, 31:3780–3792, 2022. 3

[17] Tsung-Yi Lin, Yin Cui, Serge Belongie, and James Hays. Learning deep representations for ground-to-aerial geolocalization. In *CVPR*, pages 5007–5015, 2015. 3, 4

[18] Liu Liu and Hongdong Li. Lending orientation to neural networks for cross-view geo-localization. In *CVPR*, pages 5624–5633, 2019. 3, 4

[19] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, pages 10012–10022, 2021. 3

[20] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. In *CVPR*, pages 3202–3211, 2022. 3

[21] Li Mi, Chang Xu, Javiera Castillo-Navarro, Syrielle Montariol, Wen Yang, Antoine Bosselut, and Devis Tuia. Congeo: Robust cross-view geo-localization across ground view variations. *ECCV*, 2024. 3

[22] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, pages 405–421, 2020. 3

[23] Onisimo Mutanga and Lalit Kumar. Google earth engine applications, 2019. 4

[24] Simone Alberto Peirone, Francesca Pistilli, Antonio Alliegro, and Giuseppe Averta. A backpack full of skills: Egocentric video understanding with diverse task perspectives. In *CVPR*, pages 18275–18285, 2024. 3

[25] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763, 2021. 5

[26] Krishna Regmi and Mubarak Shah. Bridging the domain gap for ground-to-aerial image matching. In *ICCV*, pages 470–479, 2019. 3

[27] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the masses. *arXiv:2104.10972*, 2021. 8

[28] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022. 5

[29] Xiaolong Shen, Zhedong Zheng, and Yi Yang. Stepnet: Spatial-temporal part-aware network for isolated sign language recognition. *IEEE TMM*, 20(7):1–19, 2024. 3

[30] Yujiao Shi, Liu Liu, Xin Yu, and Hongdong Li. Spatial-aware feature aggregation for image based cross-view geo-localization. *NeurIPS*, 32, 2019. 2, 3

[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 *CVPR*, pages 4064–4072, 2020. 3

[32] Yujiao Shi, Xin Yu, Liu Liu, Tong Zhang, and Hongdong Li. Optimal feature transport for cross-view image geo-localization. In *AAAI*, pages 11990–11997, 2020. 3

[33] Ashish Shrivastava, Tomas Pfister, Oncel Tuzel, Joshua Susskind, Wenda Wang, and Russell Webb. Learning from simulated and unsupervised images through adversarial training. In *CVPR*, pages 2107–2116, 2017. 6

[34] Jaewon Son, Jaehun Park, and Kwangsu Kim. Csta: Cnn-based spatiotemporal attention for video summarization. In *CVPR*, pages 18847–18856, 2024. 3

[35] Ze Song, Xudong Kang, Xiaohui Wei, Shutao Li, and Haibo Liu. Unified and real-time image geo-localization via fine-grained overlap estimation. *IEEE TIP*, 2024. 3

[36] Bin Sun, Chen Chen, Yingying Zhu, and Jianmin Jiang. GEOCAPSNET: ground to aerial view image geo-localization using capsule network. In *ICME*, pages 742–747, 2019. 3

[37] Yicong Tian, Chen Chen, and Mubarak Shah. Cross-view image matching for geo-localization in urban environments. In *CVPR*, pages 3608–3616, 2017. 3, 4

[38] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. *NeurIPS*, 35:10078–10093, 2022. 3

[39] Shimon Ullman. The interpretation of structure from motion. *Biological Sciences*, 203(1153):405–426, 1979. 4

[40] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017. 3

[41] Nam N Vo and James Hays. Localizing and orienting street views using overhead imagery. In *ECCV*, pages 494–509, 2016. 3, 4

[42] Tingyu Wang, Zhedong Zheng, Chenggang Yan, Jiyong Zhang, Yaoqi Sun, Bolun Zheng, and Yi Yang. Each part matters: Local patterns facilitate cross-view geo-localization. *IEEE TCSVT*, 32(2):867–879, 2021. 1, 2, 6, 7

[43] Tingyu Wang, Zhedong Zheng, Zunjie Zhu, Yuhan Gao, Yi Yang, and Chenggang Yan. Learning cross-view geo-localization embeddings via dynamic weighted decorrelation regularization. *arXiv:2211.05296*, 2022. 1, 2, 7

[44] Tingyu Wang, Zhedong Zheng, Yaoqi Sun, Chenggang Yan, Yi Yang, and Tat-Seng Chua. Multiple-environment self-adaptive network for aerial-view geo-localization. *Pattern Recognition*, 152:110363, 2024. 3

[45] Xiaolong Wang, Runsen Xu, Zhuofan Cui, Zeyu Wan, and Yu Zhang. Fine-grained cross-view geo-localization using a correlation-aware homography estimator. *NeurIPS*, 36, 2024. 2, 3

[46] Yuntao Wang, Jinpu Zhang, Ruonan Wei, Wenbo Gao, and Yuehuan Wang. Mfrgn: Multi-scale feature representation generalization network for ground-to-aerial geo-localization. In *ACM MM*, pages 2574–2583, 2024. 3

[47] Scott Workman, Richard Souvenir, and Nathan Jacobs. Wide-area image geolocalization with aerial reference imagery. In *ICCV*, pages 3961–3969, 2015. 3, 4

[48] Qiong Wu, Yi Wan, Zhi Zheng, Yongjun Zhang, Guang-shuai Wang, and Zhenyang Zhao. Camp: A cross-view geo-localization method using contrastive attributes mining and position-aware partitioning. *IEEE Transactions on Geoscience and Remote Sensing*, 2024. 3

[49] Zelong Zeng, Zheng Wang, Fan Yang, and Shin'ichi Satoh. Geo-localization via ground-to-satellite cross-view image retrieval. *IEEE TMM*, 25:2176–2188, 2022. 3

[50] Xiaohan Zhang, Xingyu Li, Waqas Sultani, Yi Zhou, and Safwan Wshah. Cross-view geo-localization via learning disentangled geometric layout correspondence. In *AAAI*, pages 3480–3488, 2023. 3

[51] Xiaohan Zhang, Xingyu Li, Waqas Sultani, Chen Chen, and Safwan Wshah. Geodtr+: Toward generic cross-view geo-localization via geometric disentanglement. *IEEE TPAMI*, 2024. 3

[52] Zhedong Zheng, Liang Zheng, and Yi Yang. Unlabeled samples generated by gan improve the person re-identification baseline in vitro. In *Proceedings of the IEEE international conference on computer vision*, pages 3754–3762, 2017. 6

[53] Zhedong Zheng, Xiaodong Yang, Zhiding Yu, Liang Zheng, Yi Yang, and Jan Kautz. Joint discriminative and generative learning for person re-identification. In *CVPR*, pages 2138–2147, 2019. 6

[54] Zhedong Zheng, Yunchao Wei, and Yi Yang. University-1652: A multi-view multi-source benchmark for drone-based geo-localization. In *ACM MM*, pages 1395–1403, 2020. 1, 3, 4

[55] Zhedong Zheng, Liang Zheng, Michael Garrett, Yi Yang, Mingliang Xu, and Yi-Dong Shen. Dual-path convolutional image-text embeddings with instance loss. *IEEE TMM*, 16 (2):1–23, 2020. 6

[56] Sijie Zhu, Taojiannan Yang, and Chen Chen. Vigor: Cross-view image geo-localization beyond one-to-one retrieval. In *CVPR*, pages 3640–3649, 2021. 3, 4

[57] Sijie Zhu, Mubarak Shah, and Chen Chen. Transgeo: Transformer is all you need for cross-view image geo-localization. In *CVPR*, pages 1162–1171, 2022. 3

# Video2BEV: Transforming Drone Videos to BEVs
# for Video-based Geo-localization

## Supplementary Material

**Outline.** This supplementary material includes two aspects:
1. Visualizations:
   - more visualizations of the Video2BEV transformation:
     - comparisons of the Video2BEV transformation at different elevation angles;
     - visualizations of drone-view videos, BEVs, and satellite-view images.
   - more visualizations of the UniV dataset;
   - more visualizations of synthetic negative samples.
2. Failure case analysis.

## 7. Visualizations

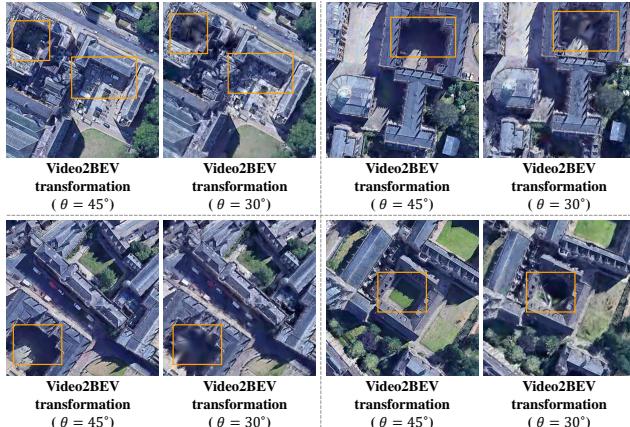### 7.1. Visualizations of the Video2BEV Transformation



Figure 8. The transformed BEV comparison of videos with different evaluation $\theta$. We highlight the challenging regions.

**Visualizations of the Video2BEV transformation at different elevation angles.** Compared to the $45°$ subset, the $30°$ subset of the UniV dataset presents more occluded cases. We analyze the impact of occlusions and other environmental constraints on the proposed Video2BEV transformation. As shown in Fig. 8, the proposed Video2BEV transformation produces satisfactory BEVs at a $45°$ elevation angle, especially in areas between tall buildings. At the $30°$ elevation angle, some regions reconstructed by the Video2BEV transformation exhibit imperfections. These imperfect regions are primarily located between buildings, where it is challenging for drones to capture clear images at a relatively low elevation angle. Despite the imperfectly reconstructed regions, the proposed Video2BEV transforma-
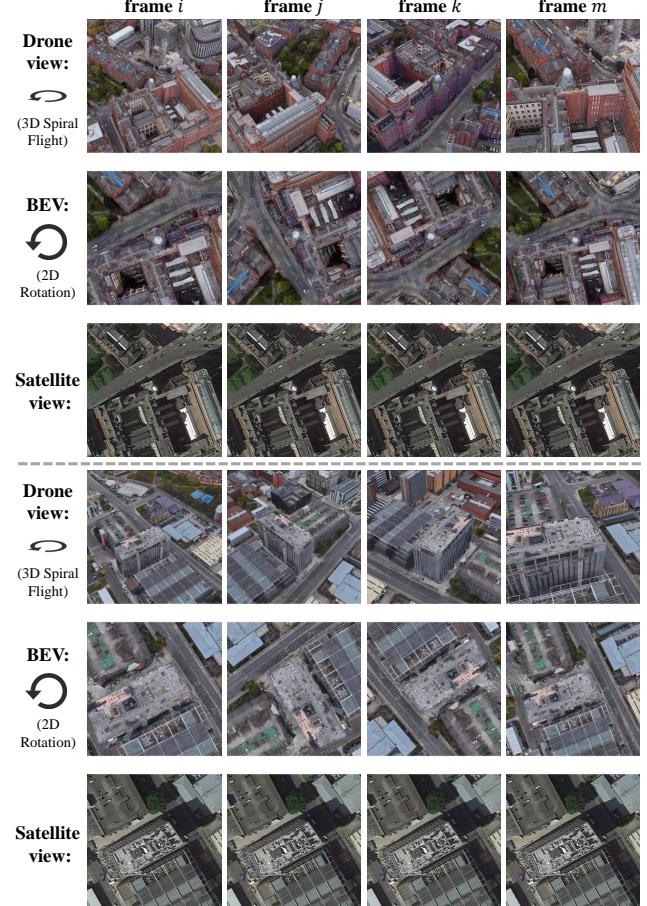


Figure 9. Visualizations of drone-view videos, BEV videos, and satellite-view images. $i, j, k, m$ are frame indices, and $i < j < k < m$.

tion significantly narrows the disparity between the drone view and the satellite view.

**Visualizations of Drone-view Videos, BEVs, and Satellite-view Images.** We provide visualizations of images from different platforms. For each building, both drone-view and Bird's Eye View (BEV) data are in video format and satellite-view data is in image format. Specifically, drones follow a spiral path around the target building, completing three circular flights. For BEVs, in the training set, we incorporate rotation angles and varying heights into BEV camera poses, generating a sequence of rotating and scaled-down BEVs (see Fig. 9). In the test set, we only increase the height of the BEV camera poses and render a

Figure 10. Elevation angles $\theta$ illustration in the UniV dataset. For each cases, top row shows $\theta = 45°$ and bottom row displays $\theta = 30°$. With a lower elevation angle, the new flight captures the target building with *wider Field of View (FoV)* but more *occlusions*, thereby posing more challenges for drone geo-localization.



Figure 11. Visualizations of original images and synthetic hard negative samples.

sequence of scaled-down BEV images. The satellite view contains one image for each building. After the proposed Video2BEV transformation, the BEVs align with the same viewing direction as the satellite view and exhibit a similar color pattern to the drone view.

### 7.2. Visualizations of the UniV Dataset

We provide additional visualizations of $45°$ and $30°$ elevation angles in the UniV dataset (see Fig. 10). Although both videos capture the same building, they differ significantly between the two elevation angles. Videos captured at a $45°$ elevation angle provide overall views of the core areas of the target building, with these areas visible in most cases. In contrast, at the relatively lower elevation angle of $30°$, drone-view videos offer a wider field of view but also introduce more occlusions. Consequently, core areas of the target building are occluded in some frames while remaining visible in others (see Fig. 10), effectively simulating outputs

from real-world drone flights.

### 7.3. Visualizations of the Synthetic Negative Samples

We provide additional visualizations of original and synthetic images (see Fig. 11). Synthetic negative samples exhibit similarities to the original samples in terms of the architectural features and color patterns of the buildings. For cases in the first row, the synthetic samples have architectural features resembling the original images, such as the circular lawn, the green sports field, and the oval stadium. In the second row, the synthetic samples exhibit similar color patterns to those of the original images. Despite the similarities, the architectural details differ between the original and synthetic samples, making the synthetic samples suitable for serving as negative samples.

### 8. Failure Case Analysis

We provide additional qualitative visualizations of retrieval results, with a particular focus on the failure cases (see Fig. 12). In these cases, the proposed method fails to recall the matched image in top-1. We observe that it is challenging because the recalled top-1 image has a very similar pattern to the query image, particularly in terms of the appearance and structure of the geographic target in the two images. In the first case, all recalled images share a similar structure and the predicted scores are relatively high. In the second case, all recalled images have a white rectangular roof. The roof of the ground truth image turns partly grey, which affects the retrieval prediction of our method. In

| BEVs via Drone Video | Satellite (Recall@1 → Recall@5) | | | | |
|---|---|---|---|---|---|
| | 0.9890 | *0.7114* | 0.5327 | 0.4813 | 0.2542 |
| | 0.0620 | *0.0538* | 0.0033 | 0.0023 | 0.0009 |

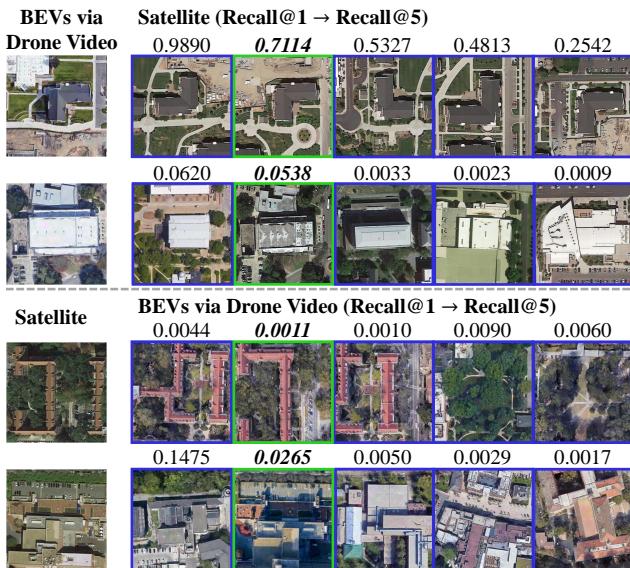| Satellite | BEVs via Drone Video (Recall@1 → Recall@5) | | | | |
|---|---|---|---|---|---|
| | 0.0044 | *0.0011* | 0.0010 | 0.0090 | 0.0060 |
| | 0.1475 | *0.0265* | 0.0050 | 0.0029 | 0.0017 |

Figure 12. Typical failure cases for Drone → Satellite and Satellite → Drone. We observe that the failures are mainly due to two factors. First, some buildings were under construction, which is quite different from the current view. Second, some satellite-view photo color is not accurate, and some similar buildings are false-matched. Given queries (left) from different platforms, matched galleries are in green box, and mismatched galleries are in blue box. The scores on the top are similarity scores estimated by the proposed method.

the third case, recalled top-3 results have similar red roofs, making it challenging for the proposed method to accurately retrieve the ground truth building. In the last case, the recalled top-1 image has a similar architectural style to the query and the ground truth image is in shadow. Both factors contribute to an inaccurate retrieval result.