

FoundLoc: Vision-based Onboard Aerial Localization in the Wild

Yao He*, Ivan Cisneros*, Nikhil Keetha, Jay Patrikar, Zelin Ye,
Ian Higgins, Yaoyu Hu, Parv Kapoor, and Sebastian Scherer

Carnegie Mellon University

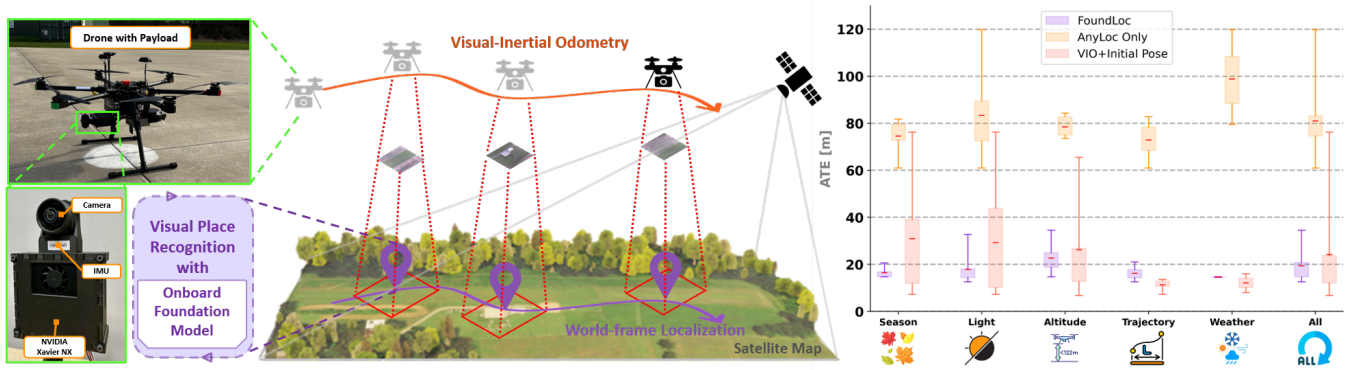


Fig. 1. *FoundLoc* enables Unmanned Aerial Vehicle global localization in the wild using only a low-cost onboard vision-based system without relying on external GNSS signals. It achieves GNSS-denied localization by anchoring a Visual-Inertial Odometry trajectory into the world frame using a Visual Place Recognition module and a satellite map (*middle*). Notably, *FoundLoc* is powered by a *foundation model* to tackle this challenging *Kidnapped Robot Problem*. In this work, we deploy our pipeline onto a custom compute-limited hardware payload and demonstrate it running in real-time (*left*). As shown in the box plot (*right*), our experiments demonstrate *FoundLoc*'s capability of dealing with the visual challenges arising from different seasons, lighting conditions, altitudes, trajectory patterns, and weather.

Abstract—Robust and accurate localization for Unmanned Aerial Vehicles (UAVs) is an essential capability to achieve autonomous, long-range flights. Current methods either rely heavily on GNSS, face limitations in visual-based localization due to appearance variances and stylistic dissimilarities between camera and reference imagery, or operate under the assumption of a known initial pose. In this paper, we developed a GNSS-denied localization approach for UAVs that harnesses both Visual-Inertial Odometry (VIO) and Visual Place Recognition (VPR) using a foundation model. This paper presents a novel vision-based pipeline that works exclusively with a nadir-facing camera, an Inertial Measurement Unit (IMU), and pre-existing satellite imagery for robust, accurate localization in varied environments and conditions. Our system demonstrated average localization accuracy within a 20-meter range, with a minimum error below 1 meter, under real-world conditions marked by drastic changes in environmental appearance and with no assumption of the vehicle's initial pose. The method is proven to be effective and robust, addressing the crucial need for reliable UAV localization in GNSS-denied environments, while also being computationally efficient enough to be deployed on resource-constrained platforms.

I. INTRODUCTION

Unmanned aerial vehicles (UAVs) have increasingly become integral in a variety of applications, ranging from agriculture to emergency response. A foundational element enabling their versatility is accurate localization, typically provided by Global Navigation Satellite Systems (GNSS). However, GNSS solutions are not without their limitations; they are vulnerable to jamming, spoofing, and environmental interference that obstruct radio signals. Therefore, the ability

of UAVs to infer their location in GNSS-denied situations is essential to creating reliable and fully autonomous systems.

Vision-based methods for UAV GNSS-denied localization are a promising solution because cameras, being passive sensors, do not suffer from the same drawbacks as GNSS-based systems, while also being low SWaP-C (Size, Weight, Power, and Cost). Extensive research in visual-inertial odometry (VIO) and Simultaneous Localization and Mapping (SLAM) [1]–[11] has yielded compelling evidence regarding the capacity of robots to achieve self-localization using only cameras and inertial measurement units (IMUs) in GNSS-denied environments. However, SLAM cannot provide Earth-fixed coordinates without external georeferencing. The accumulation of odometry drift is also a prominent concern as the robot's trajectory extends a long distance. One promising vision-based method for aerial GNSS-denied localization is Visual Terrain-Relative Navigation (VTRN) [12]–[18]. However, most VTRN methods assume the availability of initial position and heading information, as well as minimal odometric drift. This assumption permits pose propagation, and image registration or similarity estimation locally around the current position without resorting to a comprehensive global database search. Furthermore, challenges such as *appearance variance*, *stylistic dissimilarities*, and *recurring visual patterns* between vehicle camera and satellite images lead to unreliable image registration and similarity estimation. There is a need for a vision-based solution that is generalizable enough to address the above problems.

In this paper, we investigate and implement a GNSS-denied localization pipeline that relies on VIO and Vi-

*Equal Contribution

sual Place Recognition (VPR) with the generalizability of Foundation Models [19]. VPR has shown the capability of providing georeferenced image matches under large visual and viewpoint differences [20]. Unlike image registration and similarity estimation that rely on local feature retrieval and matching, VPR works with aggregated global features which are better for capturing the high-level information of an image. However, traditional learning-based VPR methods, which are usually trained on limited datasets, cannot guarantee accurate localization because of the stylistic diversity among reference satellite imagery. Furthermore, the stylistic dissimilarities between camera images and satellite images significantly reduce the matching accuracy.

Our approach, termed *FoundLoc*, is a GNSS-denied localization pipeline that attempts to tackle the aforementioned problems. The main contributions of this work are:

- We propose a novel vision-based approach to achieve GNSS-denied localization for UAVs using pre-existing satellite imagery. We design our pipeline with an In-the-Wild assumption, i.e., the vehicle does not know its initial position, and the only sensors and knowledge it has are camera images, IMU readings, and a preloaded georeferenced satellite image database.
- We formulate a Selective Ordered Top- N Recall evaluation metric, that takes into consideration whether a sufficient number of ground truth matches are present in the top N of the VPR retrieval sequence. This scoring metric helps to better evaluate VPR performance for the purposes of localization.
- We conduct extensive real-world experiments to demonstrate the robustness of our algorithm to effectively operate under In-the-Wild scenarios. With our experiments, we demonstrate the application of a foundation model in robotics tasks on computationally limited hardware. We also release our dataset, which we term “*Nardo-Air*”, that was used for these experiments.

By leveraging the strengths of a highly generalizable VPR module, our *FoundLoc* pipeline aims to address the limitations and challenges inherent in existing GNSS-denied, vision-based localization methods for UAVs, thereby offering a more robust and adaptable solution for real-world applications.

II. RELATED WORK

A. Visual-Inertial-Odometry (VIO) & VI-SLAM

Recent advancements in VIO or VI-SLAM have demonstrated remarkable performance in robot state estimation in GNSS-Denied environments. These achievements primarily leverage only camera and IMU inputs. Broadly categorized, VIO can be classified into two frameworks: filter-based [1]–[6] and graph-based optimization frameworks [7]–[11].

Filter-based approaches rely on the Extended Kalman Filter (EKF) and exhibit notable processing speed advantages over graph-based approaches. On the other hand, graph-based optimization approaches construct a factor graph of states within a given time window and perform a maximum

likelihood estimation. For nonlinear problems, graph-based approaches are more accurate than filter-based methods but require more computational resources. Both of these approaches exhibit drift due to linearization in marginalization steps and outliers [21]. Also, VIO constructs its frame with its initial state as world origin, which is independent of the Earth frame. Therefore, UAVs cannot use only VIO or VI-SLAM to obtain accurate Earth-fix coordinates.

B. Visual Terrain Relative Navigation (VTRN)

Previous vision-based UAV localization methods are usually referred to as VTRN [12]–[18]. Most VTRN methods adopt a filter-based pipeline, where odometry measurements provide information for pose updating. In pose correction, some approaches construct correctors through image registration between camera and satellite images. [12] trains a Siamese U-Net using diverse seasonal imagery datasets so that it transforms images into a season-invariant domain and conducts image registration in this domain. [13] directly conducts image registration between a camera and a satellite image by minimizing the Normalized Information Distance (NID) of the two images. [14] constructs a building ratio map (BRM) that captures the geometric features of the building area. Instead of estimating the relative transformation through registration, [14] uses the global pose estimated by BRM as a correction term.

In particle-filter-based methods, the particles are interpreted as weighted reference images or poses. [15] learns season-invariant features to calculate the similarity scores of each reference image. The similarity scores are later converted to particle weights. [16] uses normalized cross-correlation metric (NCC) to calculate image similarity and converts the similarity to particles’ likelihood. [17] generates ortho-projection from camera images and uses cross-correlation-based methods for matching score estimation between ortho-projection and satellite images. [18] estimates the Histograms of Oriented Gradients (HOG) of images and conducts a coarse to fine search for local particle weighting.

Although these VTRN methods demonstrate promising localization accuracy within their designated testing scenarios, their applicability is constrained by certain limitations that hinder their generalization. Firstly, these methods usually exhibit hard assumptions on odometry to relax the problem: (1). The initial position is available within a small error range, (2). odometry heading is known. These assumptions allow the filter to directly propagate position without a global search procedure, and image registration or particle weighting can be conducted locally around the current location. However, in terms of localization, robots are dealing with a kidnap problem where there is no prior knowledge of the initial position and heading.

On the other hand, several challenges arise due to the disparity between the nadir-facing camera images and the satellite images. These challenges can be categorized as *appearance variances*, *stylistic dissimilarities*, and *recurring visual patterns*. We will define these terminologies in Section III-E. VTRN methods are insufficient in addressing these

challenges, as they focus on local features that inadequately represent the entirety of the images.

C. Visual Place Recognition (VPR)

VPR is usually defined as an image retrieval problem where the context is to recognize previously seen places solely based on images. A basic VPR pipeline achieves image matching by computing image-wise descriptors and then calculating descriptor similarity between queries and references [22]. This pipeline offers the best trade-off between matching accuracy and search efficiency [23].

Recent VPR methods involve using a feature extraction backbone network followed by a trainable aggregation layer. One notable aggregation method is NetVLAD [24], which is a learning-based variant of the Vector-of-Locally-Aggregated Descriptors (VLAD) [25], where local features are softly assigned to a learned set of clusters. Powered by deep learning and large-scale VPR-specific data, VPR achieves a substantial boost in performance. DeLF [26] and DeLG [27] achieve large-scale image retrieval after training on Google-Landmark V1 (1 million images) and V2 datasets [28] (5 million images). TransVPR [29] and R2former [30], based on vision Transformer [31], achieve significant improvements in urban and suburban environments with 1.6 million street images from MSLS datasets [32]. Similarly, CosPlace [33] and MixVPR [34] achieve SOTA performance after training on 40 million images and 530,000 images, respectively.

Scaling up VPR-specific training datasets has been shown to be effective in improving performance. However, the aforementioned methods are usually environment-specific and task-specific, limiting their generalization capability. Trained on large and diverse datasets with self-supervision, foundation models have shown the ability to produce generalizable solutions for individual machine learning problems [35]. Benefits from the emerging foundation models, [20] proposes the first universal VPR solution AnyLoc that exhibits anywhere, anytime, and anyview capacities without any task or condition-specific training. In this paper, we adopt AnyLoc to achieve In-the-Wild satellite image recognition and provide geo-reference positions.

III. APPROACH

In our formulation, we estimate the Earth-fixed coordinates (Easting, Northing) of a UAV's position using a camera, an IMU, and pre-loaded satellite imagery. We operate under the following assumptions:

- The UAV has no knowledge of its initial position and global heading.
- Images in the satellite imagery database only provide corresponding Earth-fixed coordinates.
- The altitudes of reference images are unknown.
- The camera parameters of satellite images are unknown. We do not assume a virtual camera for satellite maps because of the stylistic difference between real camera images and satellite images.

We denote the VIO body position in the odometry frame as $\mathbf{P}_i^L = (x_i^L, y_i^L)^T$, where i indicates the i -th camera query

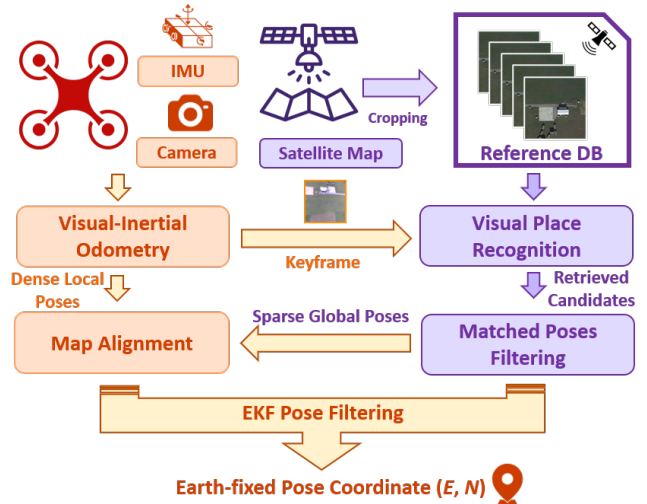


Fig. 2. The system diagram for our GNSS-Denied localization pipeline.

image and L represents the local VIO frame. The global earth-fixed coordinate provided by each satellite image is denoted as $\mathbf{P}_k^W = (x_k^W, y_k^W)$, where k indicates the k -th reference satellite image and W represents the Earth frame. $\mathcal{L}^q = \{\mathbf{I}_1^q, \mathbf{I}_2^q, \dots, \mathbf{I}_n^q\}$ denotes the set of query images and $\mathcal{L}^r = \{\mathbf{I}_1^r, \mathbf{I}_2^r, \dots, \mathbf{I}_n^r\}$ denotes the set of reference images. We denote the 6DoF $SE(3)$ pose as

$$\mathbf{T} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0}^T & 1 \end{bmatrix} \quad (1)$$

where $\mathbf{R} \in SO(3)$ is the rotation matrix and $\mathbf{t} \in \mathbb{R}^3$ is the translation in Cartesian vector space.

A. System Overview

The structure of the proposed GNSS-Denied localization pipeline is shown in Fig. 2. The VIO estimates the ego-motion of the UAV and transmits a keyframe to the VPR module at a certain frequency. The VPR module matches the query keyframe within an offline processed database where each image is cropped from the satellite map. With sufficient matches across a sequence of keyframes, the map alignment thread anchors the VIO trajectory into the Earth frame. The most recent anchored position (representing long-term memory) and the latest matching (representing instant observation) are passed to a filter to determine the UAV global position.

B. Map Alignment

The map alignment module aims to anchor odometry trajectory on the world map given a set of query images \mathcal{L}^q with their respective positions $\{\mathbf{P}_i^L\}$ in odometry frame, and a corresponding set of reference images \mathcal{L}^r with earth-fixed coordinates $\{\mathbf{P}_i^W\}$ in world frame. The reference poses are noisily corrupted, where the noise comes from:

- 1) *Cropping noise*: This noise arises from the discrete cropping of images from the satellite map, by which images are taken at certain distances. Therefore, the

observed geo-references of a query image are close to the true location, but not precisely aligned.

- 2) *False positive*: While our VPR provides SOTA recall in general, there can still be false positives in challenging areas, e.g., places that have repeated appearances.

Mathematically, the map alignment module estimates the rigid transformation T that aligns the odometry frame positions I_n^q with their corresponding noisy geo-referenced positions I_n^r . This can be expressed as

$$\min_{\mathbf{R}, t} \sum_{i=1}^N \|\mathbf{R}\mathbf{P}_i^L + t - \mathbf{P}_i^W\|_2^2 \quad (2)$$

This problem can be solved by the Iterative Closest Point (ICP) algorithm [36]. However, a common issue arises due to the nature of UAVs' predominantly movement in straight lines, causing colinearity of the query and reference image positions. This colinearity results in a degenerate scenario for the ICP algorithm, leading to rotation ambiguity.

To resolve the rotation ambiguity, we incorporate an additional constraint into the alignment problem. Specifically, we utilize gravity observability from the IMU, i.e., IMU can capture the gravity vector both in odometry and world frame. By incorporating a gravity constraint into the problem, we enforce the alignment to maintain coplanarity between the UAV's odometry trajectory and the geo-referenced map.

$$\min_{\mathbf{R}, t} \sum_{i=1}^N \|\mathbf{R}\mathbf{P}_i^L + t - \mathbf{P}_i^W\|_2^2 + \|\mathbf{R}\mathbf{g}^L - \mathbf{g}^W\|_2^2 \quad (3)$$

C. Dealing with Drift

To deal with long-term drift coming from odometry, we maintain a sliding window on the keyframe images, i.e., \mathcal{L}^q only contains at most N latest keyframes. This approach effectively reduces the impact of VIO drift over time and ensures that the query set remains relevant and up-to-date with the most pertinent keyframes.

Lastly, we conduct an EKF to update the position, where anchored odometry provides the frame-to-frame propagation. The corrector consists of two terms:

- 1) *Long-term Pose Memory*: the latest anchored position from map alignment which captures the observation distributions within the sliding window.
- 2) *Instant Pose Observation*: the latest VPR observation.

The output of EKF is the final estimated position.

D. Visual-Inertial Odometry

We use VIO to estimate the UAV ego-motion and obtain a trajectory in the odometry frame. In this process, we utilize the method proposed in [37] with a monocular camera and IMU configuration. For each image, we extract Harris corner features [38] and perform feature association using Lucas-Kanade optical flow [39] between consecutive frames. To improve robustness, the VIO incorporates outlier rejection using 5-point RANSAC [40]. Additionally, we leverage the GPU-accelerated feature extraction and tracking provided by the NVIDIA Vision Programming Interface (VPI) [41].

The pipeline described in [37] employs a sliding-window factor graph optimization, which jointly optimizes the reprojection residual, IMU preintegration residual, and a marginalization prior residual. During our experiments, we observed that optimization without a well-initialized seed can lead to frequent VIO failures, particularly during the fast flight of UAVs with computation constraints hardware. To address this issue, we perform an offline calibration of the IMU to obtain an initial estimation of the IMU acceleration bias. We define the IMU prior residual as described in Eq. (4) and incorporate it into the factor graph optimization.

$$\min_{\mathcal{T}_t} \left\{ \|\mathbf{r}_{\mathbf{b}_a}\|_{\Sigma_{\mathbf{b}_a}}^2 + \|\mathbf{r}_0\|_{\Sigma_0}^2 + \sum_{i \in \mathcal{L}} \rho(\|\mathbf{r}_{v_i}\|_{\Sigma_{v_i}}^2) + \sum_{k \in \mathcal{J}} \|\mathbf{r}_{k_{ij}}\|_{\Sigma_{k_{ij}}}^2 \right\} \quad (4)$$

$$\mathbf{r}_{\mathbf{b}_a} = \mathbf{b}_a - \hat{\mathbf{b}}_a$$

In Eq. (4), \mathcal{T}_t is the set of 6-DOF poses in the sliding window for state at current time t . \mathcal{L} and \mathcal{K} represent a set that contains visual landmarks and IMU measurements. $\rho(\cdot)$ is the robust huber loss. \mathbf{r}_{v_i} is the landmark reprojection residual, $\mathbf{r}_{k_{ij}}$ is IMU preintegration residual and \mathbf{r}_0 is the marginalization factor. Their definition can be found in [7].

$\mathbf{r}_{\mathbf{b}_a}$ is the IMU prior residual with \mathbf{b}_a and $\hat{\mathbf{b}}_a$ as the IMU acceleration bias being estimated and pre-calibrated, respectively. Enforcing a prior knowledge of IMU bias reduces VIO failures during flight. The VIO sends keyframes to the VPR at a certain rate to obtain geo-referenced images.

E. Visual Place Recognition with Foundation Model

Upon receiving a query image from VIO, we use a VPR module to obtain corresponding geo-referenced images from the onboard satellite imagery database. Challenges arise due to the disparity between the camera images and the satellite images. These challenges can be categorized as

- *Appearance variance*: The appearance of a geographical area may vary over time due to various factors such as seasonal variations, alterations in lighting conditions, changes in objects and structures, etc.
- *Stylistic difference*: Satellite maps and UAV camera images come from distinct sensing domains, which result from varying camera parameters and rendering effects. As a consequence, the visual styles differ significantly.
- *Recurring visual patterns*: There are areas that exhibit similar visual patterns such as forests, grassland, etc.

Fig. 3 demonstrates these challenges in our test area.

To address these challenges, we employ a VPR module based on the foundation model. Specifically, we utilize AnyLoc [20] with DINO Vision Transformer (ViT) [42] for dense visual feature extraction. The VPR module mainly has three stages:

1) *Offline Processing*: We generate the descriptors for each image in our database at this stage. Specifically, we adopt the aerial vocabulary from AnyLoc. This vocabulary is generated using data from both VPAIR dataset [43] and our Nardo-Air dataset. We choose DINO ViT extractor to obtain

dense visual features from images in the aerial domain. A K-means clustering on the visual features generates N_c cluster centroids, which represent the aerial vocabulary.

With the aerial vocabulary, we calculate descriptors for each satellite image in our Nardo-Air dataset. We first extract the dense features for each reference image using DINO ViT extractor. Then, we perform VLAD aggregation on these features. VLAD assigns the features to the feature vocabulary and aggregates them into image-wise VLAD descriptors.

2) *Online Inference*: During the online inference stage, we obtain the descriptor for each query image and retrieve the best matches in the database. Upon receiving a query keyframe from VIO, We extract dense visual features using DINO ViT extractor. Then we perform VLAD aggregation on the features using the aerial vocabulary. Lastly, we estimate the similarity between the query image and reference images to obtain the top- K matches with their corresponding locations in the database.

3) *False Positive Match Filtering*: We employ a density-based clustering algorithm (DBSCAN [44]) to group the set of 2D geographic points that correspond to the retrieved matches for a given query image. We remove the points that are in single-point clusters. Then, we identify the largest cluster and return the indices of the points comprising this largest group. For localization, we then use the points from this largest cluster to calculate the reference positions $\{\mathbf{P}_i^W\}$.

F. Selective Ordered Recall Metric

When evaluating and comparing VPR methods a common metric used is Recall@N, which is defined as the proportion of queries for which the correct ground-truth match appears within the top-N retrievals. Additionally, some literature such as [45] and [46] also use mAP in cases where there are many correct positive matches.

For localization purposes, where we do not employ image registration for further refinement of the database matches, we need all of the matches in the returned list of Top N matches to be close geographically, to the actual location of the query image. It is not sufficient to know whether a true-positive is somewhere in this list, but also whether there are adequate matches in the vicinity of the query image.

Depending on the discretization and overlapping strategy used when creating the database, We consider that any given query image will capture an area that overlaps with several of our database reference images. We would then expect that the reference images with the most visual overlap with the query would have the highest similarity score, and be ranked higher in the matching sequence. Thus, we propose *Selective Ordered Recall* metric, that takes into consideration whether a sufficient number of ground truth matches, say k , are present in the top N of the retrieval sequence. We denote it as Top-k@N.

We formulate the metric as follows: Let Q and R be the set of query images and the set of reference images in our database, respectively. For each query image $q \in Q$, the top- N closest retrieved reference images are denoted as

$\text{Retrieved}_N(q)$. The N ground-truth closest reference images to q are represented as $\text{GT}_N(q)$.

The Selective Ordered Recall is defined as:

$$\text{Top-k@N} = \frac{1}{|Q|} \sum_{q \in Q} \delta(q) \quad (5)$$

where $\delta(q)$ is defined as:

$$\delta(q) = \begin{cases} 1, & \text{if } |\text{Retrieved}_N(q) \cap \text{GT}_N(q)| \geq k \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

In this metric, $\delta(q)$ outputs 1 if at least k out of the top N retrieved reference images are in $\text{GT}_N(q)$, and 0 otherwise. In our implementation, we care whether at least 3 of the ground truth matches are in the top 5 retrievals. i.e., Top-3@5. We use this metric when comparing VPR performance for our pipeline in Table I.

IV. EXPERIMENTAL SETUP

We collected the “Nardo-Air” dataset, as presented in Fig. 3, for evaluation and testing of our method. Our hardware platform is described in Section IV-A. All of the data was collected at the Nardo Flight Test Field¹. This dataset comprises UAV flight data, which includes images, IMU readings, and GPS ground truth readings. It also contains a set of reference satellite imagery used for both VPR benchmarking and database descriptor computation. We describe this dataset and its collection process in the following subsections.

A. Hardware Platform

We use a custom hardware payload for both data collection and real-world flight tests. Our payload is mounted to the undercarriage of an Aurelia X6 hexacopter, in a downward-facing configuration, with the camera z-axis pointing to the ground, and the camera x-axis pointing to the left side of the hexacopter. See Fig. 1 (left) for details.

The specific components in our payload are:

- 1) *Camera Sensor*: Sony IMX264, which is a Type 2/3 11.1mm diagonal image sensor, with 5.07MP (2464 x 2056) effective pixel resolution, and a global shutter with a max frame-rate of ~ 35 FPS.
- 2) *Camera Lens*: Commonlands CIL344, Wide-Angle 4.4mm M12 lens, with F/1.9 Resolution, with IR cut-off filter, and a 100° Horizontal Field of View.
- 3) *IMU*: Epson G365, a 6 Degree of Freedom IMU with high stability, high precision, and low drift.
- 4) *GPS Module*: RadioLink SE100, a PixHawk compatible GPS receiver with up to 50cm positioning accuracy.
- 5) *Computer*: NVIDIA Jetson Xavier NX, with 384-core NVIDIA Volta™ GPU with 48 Tensor Cores and 6-core NVIDIA Carmel ARM@v8.2 64-bit CPU.

In practice, due to preprocessing, temporal alignments, and other software/firmware overhead, our effective data capture metrics were the following:

¹Location: (40.591, -79.898)

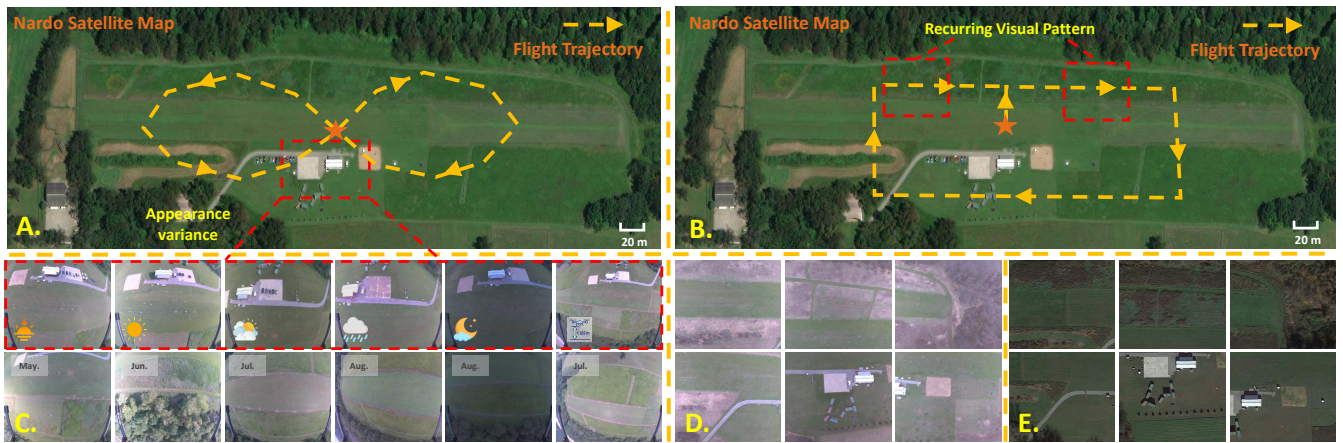


Fig. 3. **Nardo-Air** dataset. **A**: Our area of interest with one of the baseline trajectories - “Pattern Eight”. We highlight UAV imagery captured at the same location under different conditions. **B**: Our area of interest with one of the baseline trajectories - “Pattern Rectangle”. We highlight an example of a recurring visual pattern. **C**: Example camera images showcasing the different seasons, lighting variations, and altitude variations in the query images of our “Nardo-Air” dataset. **D**: The rectified query imagery (Taken in April 2023) from our test dataset used for VPR benchmarking. **E**: The cropped reference satellite imagery (Taken in November 2021 from Google Earth) from our test dataset for VPR benchmarking and online inference.

- 1) *Imagery*: 24 Hz with 1224x1028 pixel resolution.
- 2) *IMU*: Gyroscope and Accelerometer readings at 200Hz.
- 3) *GPS Coordinates*: Position readings at 1Hz.

B. Dataset Imagery

The dataset imagery consists of two sets: satellite imagery and query imagery. The sampled satellite imagery comprises the database used for image matching in the VPR module.

1) *Satellite Imagery*: Our satellite imagery is sourced from Google Maps in the form of TIF imagery and is from November 2021. We deliberately exclude the infrared channel from the images, keeping the images strictly in the RGB spectrum. This imagery boasts an interpolated spatial resolution of 0.1 m per pixel. To extract relevant samples from this source TIF, a systematic sampling strategy is employed. By overlaying a grid on our designated area of interest, we ensure structured reference TIF tiling. The spacing between each sampled image tile is 40 m in all cardinal directions: North, South, East, and West. Significantly, the central point of each image sample represents the reference position for that specific database image. Therefore, if the reference image is correct for one query, the reference position error is bounded by $40/2 = 20\text{m}$. Note that we cannot obtain altitude information in our context because satellite images usually do not provide altitude information.

Initially, every image sample extracted encompasses a dimension of 950×950 pixels. However, to streamline our analysis and maintain uniformity, these samples undergo a resizing process. Post this adjustment, each image measured 500×500 pixels. Examples of these reference tiles are shown in Fig. 3. D. In terms of physical ground coverage, each cropped and resized image spans an area with a horizontal and vertical field of view (FOV) of 60 meters. Given the set sampling distance and the FOV, it is noteworthy that there’s approximately a 33% overlap between an image and each of its immediate neighboring images. Our area-of-interest covers approximately $142,000 \text{ m}^2$.

2) *Query Imagery*: We use the onboard camera to collect another set of imagery for VPR benchmarking. We configured a lawnmower flight trajectory that covers the entire Nardo Flight Test Field at 50m altitude. The camera is north-aligned during the flight and the camera images are extracted at 1Hz. For VPR and localization metric calculation, we performed nearest neighbors assignment based on timestamps to pair GPS position messages to corresponding images.

C. UAV Trajectories & Imagery

We configure the flight trajectories and conduct aerial missions with the hexacopter utilizing a waypoint-following mode for repeatability. We record the GPS data for every trajectory as ground truth. Our experiments focus on two flight patterns: “Pattern Eight” and “Pattern Rectangle”, as illustrated in Fig. 3. A and B. Pattern Eight is 880 meters long and Pattern Rectangle is 1000 meters long. The hexacopter maintains a constant velocity of 10m/s throughout both patterns. Unless explicitly stated, the UAV maintained a consistent altitude of 50 meters above ground level during these missions.

These trajectory patterns are chosen due to their conventional but challenging shape, flight distance (which we can vary by flying multiple loops of the same path), and because they are fully contained within the safe flying zone of the Nardo Flight Test Field facility. During data collection and testing flights, we keep the yaw of the vehicle fixed, i.e., always pointing towards the North, with the assumption that the compass north-aligns images for VPR. We also keep the altitude fixed, meaning that we do not localize during the ascent/descent stages.

Our UAV dataset is comprised of 20+ trajectories with extensive flight hours from May to August 2023, covering seasonal fluctuations 🌸, diverse lighting conditions 🌞, varying altitudes \mathbb{E} (ranging from 50m to 100m), diverse flight trajectories 🛩️, and a wide range of weather conditions 🌨️, as indicated in Fig. 3. C.

V. RESULTS & DISCUSSION

We conduct a series of experiments to evaluate the accuracy and robustness of our pipeline. We first present the VPR benchmarking using both the Recall@N metric and our Selective Ordered Recall metric to demonstrate our choice of VPR methods. We further present the baseline performance comparison of different variations of our methods. We then present a comprehensive analysis of *FoundLoc* across In-the-Wild environmental conditions. Lastly, we demonstrate the real-time performance of *FoundLoc* on our hardware.

A. VPR Benchmarking

TABLE I
VPR PERFORMANCE BENCHMARK ON NARDO-AIR

Methods	R@1	R@5	Top-3@5
NetVLAD	42.25	76.06	42.25
NetVLAD-Fine-tuned	78.87	100.0	52.11
<i>AnyLoc-DINO</i>	94.37	100.0	100.0

We evaluate both Recall@N and Top-3@5 on the classical VPR method NetVLAD and AnyLoc with Dino (AnyLoc-DINO), using our Nardo-Air Dataset. As indicated in Table I, NetVLAD has a low R@1, R@5, and Top-3@5 without any data-specific fine-tuning. This is primarily attributed to the inherent challenges stemming from the dissimilarities between camera images and satellite images. With fine-tuning on our dataset, NetVLAD archives a significant improvement on R@1 and has 100% R@5. However, there is a minor improvement on Top-3@5, indicating that there is a concerning amount of false positives among the top-5 retrievals even with model fine-tuning. For localization purposes, a false positive will significantly reduce the accuracy, which is intolerable. On the other hand, AnyLoc-DINO has a significantly higher Recall@N and 100% Top-3@5 even without model fine-tuning. This guarantees reliable georeferences during the flight.

B. *FoundLoc* Performance Analysis & Baseline Comparison

In this section, we analyze *FoundLoc*'s performance across In-the-Wild conditions and compare it with baseline methods. Note that VTRN methods have strong assumptions on initial robot poses and use simulated odometry with known heading, so VTRN methods fail in our scenarios. We also evaluate the effect of altitude changes, which VTRN methods do not consider. The baseline methods include: VIO+IP (VIO with known initial pose, i.e., position and orientation), VIO+NetVLAD (VIO with Fine-tuned NetVLAD to provide georeference), AnyLoc-DINO (average positions of Top-3 image retrieves from AnyLoc with DINO ViT extractor), and *FoundLoc*-NF (*FoundLoc* with no false positive filtering). We use ATE as our comparison metric. The results are shown in Table II, with Fig. 1 (right) plotting the ATEs.

1) *Accuracy & Robustness*: *FoundLoc* achieves an average ATE below 20m at various seasons, lighting conditions, trajectory patterns, and various weather conditions, with the standard deviation of all the scenarios lower than 10m. *FoundLoc* maintains a consistent performance under these scenarios, demonstrating its accuracy and robustness under In-the-Wild conditions. Note that the ATE errors echo with the positional error bound in Section IV-B.1. The ATE at altitude changes is slightly above 20m. This is because we do not conduct drone tilting correction. As a consequence, the keyframe sent to VPR is not exactly below the drone, causing drift when obtaining reference positions.

2) *Compared with VIO only Approach*: While VIO can achieve relative accuracy with known initial poses, it exhibits higher ATE compared to *FoundLoc* across various scenarios. Furthermore, the standard deviation of ATE in VIO is notably larger than that of *FoundLoc*. This discrepancy is primarily attributed to VIO's susceptibility to drift and scale errors, even when provided with a strong initial pose assumption. VIO tends to exhibit significant drift over the course of UAV flights. Additionally, high-altitude and high-speed flight hinder accurate scale estimation. In contrast, *FoundLoc* benefits from continuous drift correction through Visual Place Recognition (VPR), which effectively mitigates drift and scale errors.

3) *Compared with VPR Methods and Filtering Effect*: Even with 78.87 R@1 on the VPR benchmark dataset, NetVLAD struggles to provide reliable geo-reference positions from the satellite map. This is mainly due to the limited generalization capability of NetVLAD. In our experiments, we observe that NetVLAD keeps retrieving reference images from the center of the satellite map (framed part of Fig. 3.A). This recurrent behavior results in failures in map alignment and, consequently, disrupts the entire pipeline.

Powered by the foundation model, AnyLoc is capable of generalizing to satellite image retrieval tasks. However, it still exhibits a large amount of false positives, decreasing the localization accuracy. This false positives effect is also reflected in *FoundLoc*-NF. With our false positive filtering method, *FoundLoc* significantly reduce the localization error.

C. Real-time Performance of *FoundLoc*

We deploy our module to onboard computational resource-limited hardware. The VIO estimates the pose at 8Hz with a total CPU consumption between 300% and 400% (600% maximum). The run-time performance of the foundation model is our interest. We use the ViT-S/8 model with 21 million parameters for dense DINO feature extraction. The input images are downsampled from 500×500 pixels to 224×298 pixels. The foundation-model-based VPR achieves 1.934Hz inference frequency with above 90% GPU usage on NVIDIA Xavier NX. VLAD encoding and matching take 0.107s and 0.0018s per query image, respectively.

VI. CONCLUSION

In this paper, we design and implement a vision-based GNSS-denied localization pipeline for UAVs that uses a

TABLE II

ATE COMPARISON ON NARDO-AIR DATASET IN METERS. THE METRICS ARE AVERAGE ATE (AVE.) AND STANDARD DEVIATION OF ATE (SD.)

Methods	Season 🌸		Light 🌅		Altitude 🏔️		Trajectory 🗺️		Weather 🌤️		All 🌐	
	AVE.	SD.	AVE.	SD.	AVE.	SD.	AVE.	SD.	AVE.	SD.	AVE.	SD.
VIO + IP	30.92	31.26	29.26	24.21	26.26	30.91	11.32	1.29	12.07	3.97	24.18	22.34
VIO + NetVLAD	-	-	-	-	-	-	-	-	-	-	-	-
AnyLoc-DINO	74.56	9.27	83.34	26.35	78.46	5.85	72.91	4.16	98.81	20.16	80.96	13.33
<i>FoundLoc</i> -NF	46.61	84.75	41.07	57.14	57.75	90.14	84.51	100	24.02	38.43	27.72	49.50
<i>FoundLoc</i>	16.41	2.82	17.86	9.50	22.69	6.25	16.20	3.68	14.63	0.14	19.38	6.11

downward-facing camera and an IMU to accurately and robustly localize in different environments. We combine a robust VIO and a robust VPR with a map alignment module to get a global estimate of our pose on a map using an onboard satellite imagery database. We utilize a foundation-model-based image matcher, that allows our system to accurately anchor with respect to *appearance variances, stylistic dissimilarities, and recurring visual patterns* without any task or condition-specific model tuning. Our tests demonstrate that our pipeline can generate accurate localization results. We also deploy our pipeline onto an onboard system and demonstrate the possibility of running a foundation model on computation-limited hardware.

ACKNOWLEDGMENTS

Approved for public release; distribution is unlimited. This research was sponsored by DARPA (W911NF-18-2-0218). The views, opinions, and/or findings expressed are those of the author(s) and should not be interpreted as representing the official views or policies of the Department of Defense or the U.S. Government. The authors thank Jay Karhade, Avneesh Mishra, Krishna Murthy Jatavallabhula, & Alaa Maalouf for their support with the deployment of AnyLoc. We also thank John Keller for helping out with the hardware.

REFERENCES

- [1] A. I. Mourikis and S. I. Roumeliotis, "A multi-state constraint kalman filter for vision-aided inertial navigation," in *Proceedings 2007 IEEE international conference on robotics and automation*. IEEE, 2007, pp. 3565–3572. 1, 2
- [2] A. J. Davison, I. D. Reid, N. D. Molton, and O. Stasse, "Monoslam: Real-time single camera slam," *IEEE transactions on pattern analysis and machine intelligence*, vol. 29, no. 6, pp. 1052–1067, 2007. 1, 2
- [3] M. Bloesch, S. Omari, M. Hutter, and R. Siegwart, "Robust visual inertial odometry using a direct ekf-based approach," in *2015 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2015, pp. 298–304. 1, 2
- [4] E. S. Jones and S. Soatto, "Visual-inertial navigation, mapping and localization: A scalable real-time causal approach," *The International Journal of Robotics Research*, vol. 30, no. 4, pp. 407–430, 2011. 1, 2
- [5] J. A. Hesch, D. G. Kottas, S. L. Bowman, and S. I. Roumeliotis, "Camera-imu-based localization: Observability analysis and consistency improvement," *The International Journal of Robotics Research*, vol. 33, no. 1, pp. 182–201, 2014. 1, 2
- [6] P. Geneva, K. Eickenhoff, W. Lee, Y. Yang, and G. Huang, "Opennvins: A research platform for visual-inertial estimation," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 4666–4672. 1, 2
- [7] T. Qin, P. Li, and S. Shen, "Vins-mono: A robust and versatile monocular visual-inertial state estimator," *IEEE Transactions on Robotics*, vol. 34, no. 4, pp. 1004–1020, 2018. 1, 2, 4
- [8] C. Campos, R. Elvira, J. J. G. Rodríguez, J. M. Montiel, and J. D. Tardós, "Orb-slam3: An accurate open-source library for visual, visual-inertial, and multimap slam," *IEEE Transactions on Robotics*, vol. 37, no. 6, pp. 1874–1890, 2021. 1, 2
- [9] A. Rosinol, M. Abate, Y. Chang, and L. Carlone, "Kimera: an open-source library for real-time metric-semantic localization and mapping," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 1689–1696. 1, 2
- [10] C. Forster, M. Pizzoli, and D. Scaramuzza, "Svo: Fast semi-direct monocular visual odometry," in *2014 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2014, pp. 15–22. 1, 2
- [11] L. Von Stumberg, V. Usenko, and D. Cremers, "Direct sparse visual-inertial odometry using dynamic marginalization," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 2510–2517. 1, 2
- [12] A. T. Fragoso, C. T. Lee, A. S. McCoy, and S.-J. Chung, "A seasonally invariant deep transform for visual terrain-relative navigation," *Science Robotics*, vol. 6, no. 55, p. eabf3320, 2021. 1, 2
- [13] B. Patel, T. D. Barfoot, and A. P. Schoellig, "Visual localization with google earth images for robust global pose estimation of uavs," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 6491–6497. 1, 2
- [14] J. Choi and H. Myung, "Brm localization: Uav localization in gnss-denied environments based on matching of numerical map and uav images," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 4537–4544. 1, 2
- [15] J. Kinnari, F. Verdoja, and V. Kyrki, "Season-invariant gnss-denied visual localization for uavs," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 10232–10239, 2022. 1, 2
- [16] R. Jurevičius, V. Marcinkevičius, and J. Šeibokas, "Robust gnss-denied localization for uav using particle filter and visual odometry," *Machine Vision and Applications*, vol. 30, pp. 1181–1190, 2019. 1, 2
- [17] J. Kinnari, F. Verdoja, and V. Kyrki, "Gnss-denied geolocalization of uavs by visual matching of onboard camera images with orthophotos," in *2021 20th International Conference on Advanced Robotics (ICAR)*. IEEE, 2021, pp. 555–562. 1, 2
- [18] M. Shan, F. Wang, F. Lin, Z. Gao, Y. Z. Tang, and B. M. Chen, "Google map aided visual navigation for uavs in gps-denied environment," in *2015 IEEE international conference on robotics and biomimetics (ROBIO)*. IEEE, 2015, pp. 114–119. 1, 2
- [19] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill *et al.*, "On the opportunities and risks of foundation models," Aug. 2021. 2
- [20] N. Keetha, A. Mishra, J. Karhade, K. M. Jatavallabhula, S. Scherer, M. Krishna, and S. Garg, "Anyloc: Towards universal visual place recognition," *arXiv preprint arXiv:2308.00688*, 2023. 2, 3, 4
- [21] C. Forster, L. Carlone, F. Dellaert, and D. Scaramuzza, "On-manifold preintegration for real-time visual-inertial odometry," *IEEE Transactions on Robotics*, vol. 33, no. 1, pp. 1–21, 2016. 2
- [22] S. Schubert, P. Neubert, S. Garg, M. Milford, and T. Fischer, "Visual place recognition: A tutorial," *arXiv preprint arXiv:2303.03281*, 2023. 3
- [23] S. Garg, T. Fischer, and M. Milford, "Where is your place, visual place recognition?" *arXiv preprint arXiv:2103.06443*, 2021. 3
- [24] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "Netvlad: Cnn architecture for weakly supervised place recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5297–5307. 3

- [25] R. Arandjelovic and A. Zisserman, "All about vlad," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2013, pp. 1578–1585. 3
- [26] H. Noh, A. Araujo, J. Sim, T. Weyand, and B. Han, "Large-scale image retrieval with attentive deep local features," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 3456–3465. 3
- [27] B. Cao, A. Araujo, and J. Sim, "Unifying deep local and global features for image search," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XX 16*. Springer, 2020, pp. 726–743. 3
- [28] T. Weyand, A. Araujo, B. Cao, and J. Sim, "Google landmarks dataset v2—a large-scale benchmark for instance-level recognition and retrieval," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2575–2584. 3
- [29] R. Wang, Y. Shen, W. Zuo, S. Zhou, and N. Zheng, "Transvpr: Transformer-based place recognition with multi-level attention aggregation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 13 648–13 657. 3
- [30] S. Zhu, L. Yang, C. Chen, M. Shah, X. Shen, and H. Wang, "R2former: Unified retrieval and reranking transformer for place recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 19 370–19 380. 3
- [31] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020. 3
- [32] F. Warburg, S. Hauberg, M. Lopez-Antequera, P. Gargallo, Y. Kuang, and J. Civera, "Mapillary street-level sequences: A dataset for lifelong place recognition," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2626–2635. 3
- [33] G. Berton, C. Masone, and B. Caputo, "Rethinking visual geolocalization for large-scale applications," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 4878–4888. 3
- [34] A. Ali-Bey, B. Chaib-Draa, and P. Giguere, "Mixvpr: Feature mixing for visual place recognition," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 2998–3007. 3
- [35] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill *et al.*, "On the opportunities and risks of foundation models," *arXiv preprint arXiv:2108.07258*, 2021. 3
- [36] K. S. Arun, T. S. Huang, and S. D. Blostein, "Least-squares fitting of two 3-d point sets," *IEEE Transactions on pattern analysis and machine intelligence*, no. 5, pp. 698–700, 1987. 4
- [37] Y. He, H. Yu, W. Yang, and S. Scherer, "Towards robust visual-inertial odometry with multiple non-overlapping monocular cameras," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2022, pp. 9452–9458. 4
- [38] C. Harris, M. Stephens *et al.*, "A combined corner and edge detector," in *Alvey vision conference*, vol. 15, no. 50. Citeseer, 1988, pp. 10–5244. 4
- [39] B. D. Lucas and T. Kanade, "An iterative image registration technique with an application to stereo vision," in *IJCAI'81: 7th international joint conference on Artificial intelligence*, vol. 2, 1981, pp. 674–679. 4
- [40] M. A. Fischler and R. C. Bolles, "Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography," *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981. 4
- [41] NVIDIA. (2023) Vpi-vision programming interface documentation. Accessed on: 2023-08-15. [Online]. Available: <https://docs.nvidia.com/vpi/index.html> 4
- [42] M. Caron, H. Touvron, I. Misra *et al.*, "Emerging properties in self-supervised vision transformers," in *ICCV*, 2021. 4
- [43] M. Schleiss, F. Rouatbi, and D. Cremers, "Vpair—aerial visual place recognition and localization in large-scale outdoor environments," *arXiv preprint arXiv:2205.11567*, 2022. 4
- [44] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise," <https://cdn.aaii.org/KDD/1996/KDD96-037.pdf>, accessed: 2023-9-11. 5
- [45] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman, "Object retrieval with large vocabularies and fast spatial matching," in *2007 IEEE Conference on Computer Vision and Pattern Recognition*, Jun. 2007, pp. 1–8. 5
- [46] H. Jegou, M. Douze, and C. Schmid, "Hamming embedding and weak geometric consistency for large scale image search," in *Computer Vision – ECCV 2008*. Springer Berlin Heidelberg, 2008, pp. 304–317. 5