

Key-Region-Based UAV Visual Navigation

Michael Karnes¹, Jacob Riffel², Alper Yilmaz³

¹Dept. of Civil Engineering, The Ohio State University, 281 W Lane Ave, Columbus, Ohio - karnes.30@osu.edu

²Dept. of Computer Science and Engineering, The Ohio State University, 281 W Lane Ave, Columbus, Ohio - riffel.8@osu.edu

³Dept. of Civil Engineering, The Ohio State University, 281 W Lane Ave, Columbus, Ohio - yilmaz.15@osu.edu

Keywords: Visual Navigation, Visual Geolocalization, Few-shot Learning Re-Identification, DNN

Abstract

Visual navigation has recently seen significant developments with the rise in autonomous navigation. Keypoint-based mapping and localization has served as a reliable localization method for many applications, but the push to run more applications on less expensive hardware becomes extremely limiting. In this paper, we present a novel approach for visual geolocalization and navigation that improves landmark detection reliability while reducing reference map complexity. Similar to prior techniques, we use the process of point based matching schemes to solve for the image-to-map transform. The critical difference is that we use object detection to identify key-regions instead of keypoints. During an initial flight key-regions are mapped into an identity dictionary with their geolocations and few-shot learning encoded descriptors. Then on subsequent flights, key-regions are detected and matched using the identity dictionary for re-identification. Using the identified vehicles as key-regions, the results show that the proposed key-region based localization produces GPS like localization while maintaining a higher resilience to image noise compared to keypoint-based techniques.

1. Introduction

Geolocalization is a critical tool in today's modern economy with a long list of direct applications and impacts in logistics, construction, analytics, and navigation. The impact of publicly available Global Positioning System (GPS) is hard to overstate, bringing global collaborators to a unified common coordinate system. As researchers have extended the utility of GPS to more precise and accessible applications, the challenges associated with GPS become apparent. This is especially true in the effort of autonomous unmanned aerial vehicles (UAV) where operating environments and precision are being pushed to the limit.

Autonomous drones have been utilized in aerial, outdoor, indoor, underground, and underwater applications. Aerial applications have the most reliable GPS reception with little interference. As the operating environment gets closer to the ground, the likelihood of obstructing the satellite signals increases, such as in urban or mountainous environments. Indoor, underground, and underwater render GPS extremely unreliable.

As the operating environments and tasks become more complex, the need for precision and reliability increases. For example, UAV package delivery in populated areas has significant risks of human injury if its localization system were to fail. As the adoption of autonomous UAVs in working environments increases, so does the need for reliable localization.

Visual localization provides an alternative approach to GPS. This is commonly accomplished with well defined landmarks, such as seen in pre-GPS aviation with human readable labeled ground control points. For computer vision, landmarks are often translated into machine interpretable fiducial markers that have been placed and mapped throughout the operating environment. This approach is suitable for many controlled environments but has significant challenges in scaling to larger more complex environments.

Simultaneous localization and mapping (SLAM) provides a particularly flexible approach toward visual navigation, producing simultaneous point cloud mapping and camera localization. This approach is effective but notoriously computationally expensive. This is partially due to the curse of dimensionality. SLAM is based on identifying keypoints that serve as naturally occurring landmarks. The task of identifying, matching, and mapping can rapidly get out of hand if not delicately balanced. The implementation of SLAM requires the user to carefully consider the scene characteristics and tune the decision hyperparameters accordingly.

Some more recent approaches improve efficiency by adapting the SLAM approach to a fixed reference map. This reduces the problem complexity and counteracts the effects of accumulating error. This approach hits a balance of performance and reliability, utilizing pre-existing georectified imagery to remove the need for manual fiducial marker installment. However, it still suffers from sensitivity to appearance variation requiring the reference imagery to be taken in similar conditions as the expected use case.

Key-region visual navigation provides the advantages of reference map based visual navigation while reducing sensitivity to appearance changes and improving the understandability of extracted landmarks. This is accomplished by extending the keypoint concept into a more stable characteristic attribute space. Keypoints identify distinct points within the scene and generate a discriminative description of the neighborhood around that keypoint for later re-identification. The small area of influence of the keypoint makes it highly sensitive to perspective and lighting changes.

An additional advantage to key-region visual navigation is its understandability. The concept of a keypoint is highly ambiguous as to what it semantically refers. Good keypoints are scene points that are the most reliably recognizable leaving the user without a clear definition of what to look for when assessing

operations in a new environment. Key-regions enable the user to associate landmarks with semantically meaningful objects, making their assessment much more reasonable. For example, if buildings are used as key-regions the user can easily assess the number and uniqueness of the buildings within a scene.

This paper presents a novel key-region based visual navigation algorithm. This algorithm leverages object detection networks to extract full objects as passive landmarks from a scene seen in Figure 1. Few-shot learning is then used to encode the detected objects into discriminative descriptors for later re-identification. This approach is agnostic and can be generalized across a wide range of objects. Detection networks have been well established at detecting a large variety of different objects, and the few-shot learning approach is specifically developed to rapidly adapt to novel concepts with a few examples images.

Cars from a UAV perspective were selected as the example key-region object for several reasons, including the availability of public datasets to train the detection network, their high frequency of occurrence across operating environments, and their high visibility within the scene. The aim of this paper is to present the generalizable key-region based visual navigation that can be adapted for a user's specific application.

Estimated GPS: (40.047141, -83.128684)



Figure 1. Example visualization of the key-region based localization.

2. Background

UAVs have a variety of localization capabilities. The standard modern UAV (DJI, Parrot) comes with onboard GPS and inertial measurement unit (IMU) based localization systems. (Jametoni and Saputra, 2021) lays out the remaining localization alternatives in a taxonomy that categorizes UAV localization systems into two major methods: vision-based and non-vision-based. The non-vision based refers to these GPS, IMU,

and radio frequency (RF) based transmission networks; while the vision based refers to ground control points, fiducial markers, and SLAM. (Czyza et al., 2023) provides an assessment of the state-of-the-art DJI Matrice 300 real-time kinematic positioning (RTK) localization validating the published specs for a GPS horizontal localization error of 0.50 m and an RTK corrected average error of 0.10 m. (Elkhrachy, 2021) backs these findings by exploring the relative localization error of the DJI Mavic Pro Platinum using 21 ground control points with root mean square error (RMSE) of 0.883 m. (Patrik et al., 2019) tested the positional accuracy of Global Navigation Satellite System (GNSS) based autonomous flight on the Erle Robotics UAV finding an average horizontal localization error of 1.11 m. In general, these onboard GNSS based systems produce a mean error of 0.5-2 m.

Many efforts are looking to improve UAV localization through the addition of control beacons. These beacons take both visual and non-visual forms. For example, the addition of ultra wide-band (UWB) beacons were shown to reduce localization errors from 0.35 m using the original GPS to 2.02 m with beacons (Chen et al., 2023). A wide variety of RF beacons are currently used but can come with increased complexities and localization uncertainties, such as seen when using Bluetooth low energy (BLE) and WiFi signals (Chen et al., 2023).

Much in the same way as RF, visual fiducial markers have been employed to reduce localization errors (Mráz et al., 2020, ?). (Supriyono and Akhara, 2021) presented a direct comparison between GPS and GPS + visually aided landing, showing a reduction in localization error of 1.99 m to 0.45 m. Beacon placement is an effective high precision localization approach. However, the installation and maintenance of beacon systems are costly and infeasible for many UAV applications. This is where the advantage of passive keypoint based visual navigation methods becomes apparent.

Keypoint based visual navigation does not require prior installations of beacons. Instead naturally occurring beacons/landmarks are extracted from the scene using a keypoint generator. These keypoints are generally selected by the patterning in a pixel neighborhood. This can be in the form of handcrafted gradient based features such as SIFT and ORB or deep neural network (DNN) based features such as SuperPoint (Wang et al., 2024, Wei and Yilmaz, 2023). These methods search the image space, identify pixels with distinct neighborhoods, and then generate a description vector of that pixel region by which it can later be re-identified.

(Gupta and Fernando, 2022) provides a review of the current state-of-the-art SLAM approaches, giving an overview of the approach and applications. These approaches vary from solely SLAM, SLAM integrated GNSS, SLAM with multi-sensor fusion, and SLAM for UAV swarms. These types of systems are capable of achieving precise localization of 0.33 m (Li et al., 2022). Another recent advancement combined SLAM and open street map building geometry (Frosi et al., 2023). Another direction improves the efficiency of SLAM by introducing a pre-flight generated keypoint reference map created from prior captured satellite imagery (Wei and Yilmaz, 2023). This approach demonstrated strong geolocalization with an accuracy of 3.4 m over km long flights in urban and rural environments. Another study combined SLAM with object recognition to map the 3D location of detected objects (Mazurek and Hachaj, 2021).

Key-region based visual navigation extends the keypoint concept

towards key-objects. This approach leverages the recent developments in visual navigation with the addition of more stable landmarks. Keypoints are relatively smaller than objects capturing distinct information directly around a pixel making them sensitive to perspective variation. Objects are recognizable from a wider range of perspectives, ranges, and viewing conditions. (Pi et al., 2020) completes the task of object re-identification through simultaneous multiple views. Through this process, they were able to accurately identify and map the scene objects but did not address UAV geolocalization over a flight. The key development that enables this approach is the advancement in rapid object re-identification providing the means for object matching.

Key-region based visual navigation map utilizes few-shot generated appearance models for object re-identification. Few-shot learning is a growing branch of the visual classification field. Originating from the desire to reduce training requirements of classification DNN, few-shot learning reformulated the task into a metric learning form. The goal changed to focus on rapid learning through optimized discriminative feature embedding (Zeng and ying Xiao, 2024). Later this changed to a meta-learning approach where the goal is to design models that can be rapidly retrained (Zeng and ying Xiao, 2024).

For this work, we selected a few-shot learning method with high generality and minimal training requirements known as the Omni-Modeler (Karnes and Yilmaz, 2023). This approach directly extracts and transforms a discriminative feature set from the generic pretrained VGG-11 model without requiring iterative training or annotated data. The feature transform process is referred to as language encoding. The key-region mapping process intakes generic object detections and stores them in an identity dictionary along with their encoded descriptions. On following flights over the area, new detections are matched by comparison to the identity dictionary.

UAV autonomous navigation heavily relies on GPS/GNSS localization with the growing recognition that higher precision and more reliable operation are essential for emerging applications. Beacon and ground control point based navigation are highly effective but dependent upon maintaining a reliable infrastructure. Keypoint based approaches relieve these infrastructure needs but are highly sensitive to scene changes limiting their operational bounds. Key-region based navigation leverages this long line of UAV navigation along with advancements in object recognition to provide a more robust and understandable visual navigation strategy.

3. Algorithm Design

3.1 Overview

Key-region based visual navigation has two primary phases: mapping and localization. The mapping phase is completed on an initial flight capturing imagery and location data on the objects of interest. This information is then processed into the identity dictionary, a JSON file describing the mapped area in terms of the objects' sizes, locations, and descriptions. This map is then re-utilized on the same UAV or transferred to an additional UAV for future navigation through the area.

3.2 Map Generation

The key-region object map is generated during a mapping flight where other localization methods are available. During this pro-

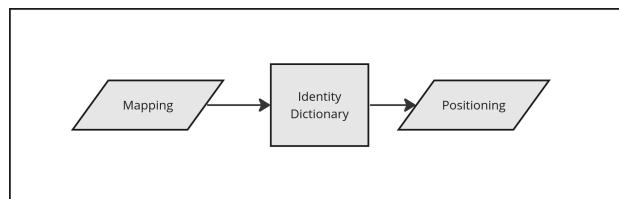


Figure 2. Algorithm overview showing the mapping to identity dictionary to positioning relationship.

cess, the UAV position is determined using a keypoint localization method similar to (Wei and Yilmaz, 2023). Then each frame is passed through the object detection network, in this case, the SAHI car detection network (Akyon et al., 2021). The bounding box for each detection is then projected into the geocoordinate space using the homography calculated from the keypoints. Each detection is stored in a detections aggregation dictionary. After the flight is completed, object detections are associated into identities using their mapped bounding box areas. Once the identities are formed, their corresponding image regions are cropped from the original image and encoded into a description. Each identity is stored in the identity dictionary along with its bounding box information, mapped bounding box information, and descriptions.

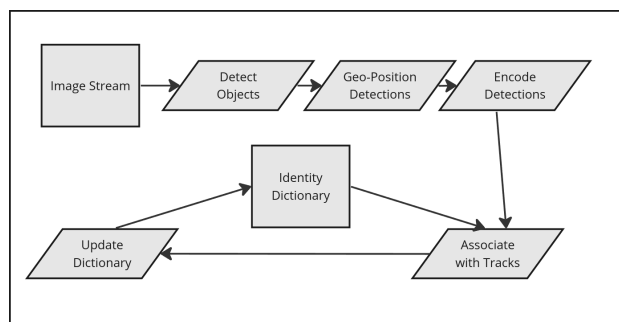


Figure 3. Flowchart showing the map generation process.

3.3 Description Encoding

Object descriptions are encoded using the Omni-Modeler (Karnes and Yilmaz, 2023). The Omni-Modeler is a generalized few-shot learning algorithm designed to rapidly adapt to novel objects. This is done through a series of calculated transforms that maximize the discriminability of the resulting feature space. This process is modeled off of the structure of natural language, moving from letters to words to sentences. Using a set of representative unannotated data, first the letters are extracted from the raw latent feature space of the general pretrained DNN using ICA. Next, the words are extracted from the letter feature space using k-means clustering. Finally, the words are combined into sentences by concatenating across DNN layers.

After the language encoding transforms are calculated, then an identity dictionary is created. This process takes in example reference images for each class and aggregates their associated detections into a dictionary. The positional data and encoded descriptions for each detection are stored for each identified object. New samples are then identified by comparison to the identity dictionary.

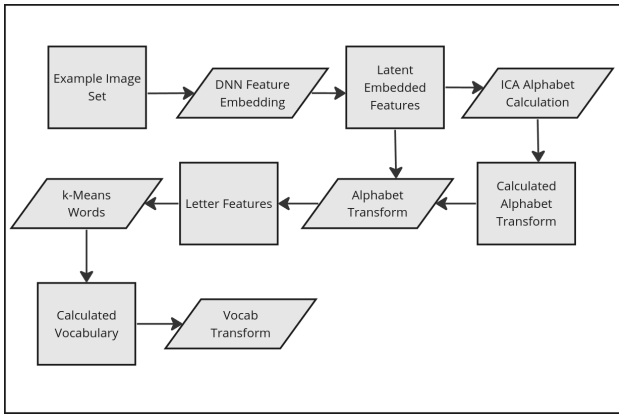


Figure 4. Flowchart showing the Omni-Modeler training.

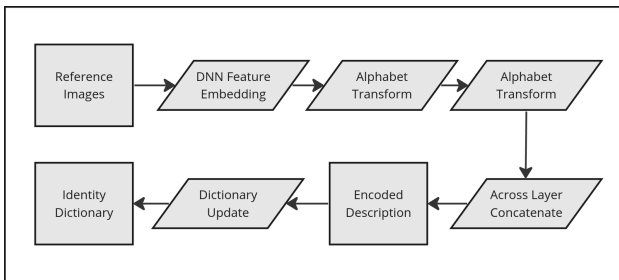


Figure 5. Flowchart showing the Omni-Modeler encoding process.

3.4 Localization

Key-region based visual navigation follows a similar process as keypoint based navigation approaches, which detects points in the UAV image stream, matches them to prior known points on the reference map, and then calculates the affine projection transform of the camera image to the reference map points. In the same manner, key-region based navigation detects key-objects, and then identifies matches and positions from those matches.

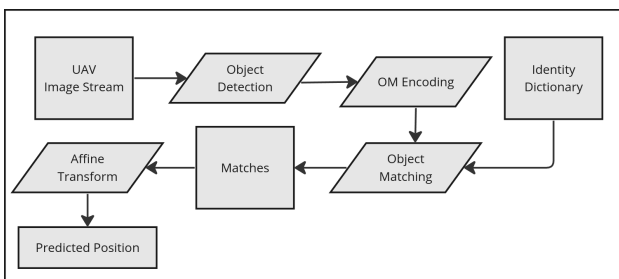


Figure 6. Flowchart showing the key-region based localization process.

4. Experiments

4.1 Overview

The key-region based navigation approach created several decision point hyperparameters for the mapping and localization processes. The experiments in this study explore the most influential decision parameters for each of these processes. The

Intersection-over-Union (IoU) Association Threshold controls the decision point at which detections are associated into identities. The Probability Matching Threshold controls the matching decision process critical to final localization performance.

The effects of these hyperparameters are explored using two flight paths over the same area flown on the same afternoon with the camera facing in a Nadir view. The first flight path, flown at 300 ft and used for the mapping phase, is referred to as Flight Path #1. The second flight path, used for the localization phase, is referred to as Flight Path #2, flown at both 200 ft and 300 ft.



Figure 7. Visualizations of the two flight-paths utilized in this study. The top image shows Flight Path #1, flown at 300 ft and used for the mapping process. The bottom image shows the flight path for Flight Path #2 which was flown at two different altitudes, 200 ft and 300 ft, for localization testing.

4.2 Mapping

The mapping process creates the identity dictionary that consists of a series of object entries with the position and description data for each of their detections. An object identity is generated from identities of associated detections. The mapping flight produces a series of detections and their mapped positions. The mapped detections are associated into identities by their IoU values. This study looks at the quantitative and qualitative effects of varying the IoU Association Threshold from 0.3 to 0.8 with a step size of 0.1.

Table 1 shows the results of the varying IoU Association Threshold, showing the number of extracted identities and the average number of detections per identity. As the threshold increases the number of identities and their average number of associated detections decreases.

IoU Thresh.	Identities	Avg. Dets.
0.3	43	141.37
0.4	42	143.48
0.5	43	136.60
0.6	40	139.68
0.7	42	126.29
0.8	37	122.41

Table 1. IoU Association Threshold.

Figure 8 shows a visualization of the extracted identities from the mapping process when using an IoU threshold of 0.1 and 0.8 where each of the blue bounding boxes represents an extracted identity. The largest qualitative difference seen between the two cases is the reduction of overlapping identity bounding boxes as well as their relative sizes and aspect ratios. Combining these quantitative and qualitative results it appears that the higher 0.8 IoU threshold produced a higher quality map.

4.3 Localization

The localization process uses the identity dictionary as its key-region reference map. As the flight progresses, the algorithm compares new detections to its reference to determine matches. The comparison is completed using the 5 - Nearest Neighbors which provides the most likely match and its associated probability. The associated probability serves as a match quality score. This localization experiment looks into the effect of varying Probability Matching Threshold.

For this study, the identity dictionary resulting from Flight Path #1 with an IoU threshold of 0.8 was used to localize the 200 ft and 300 ft flights over Flight Path #2. Table 2 shows the effects of varying the Probability Matching Threshold on localization performance. The 200 ft Flight Path #2 showed an overall weaker performance. As the threshold increases, the mean localization error also tends to increase. The minimum mean error for this flight is 3.36 m with probability thresholds of 0.3 and 0.4. The opposite trend was seen for the 300 ft Flight Path #2, with a continually decreasing mean error as the probability threshold increases. The best performance for this flight has a mean error of 2.23 m with a probability threshold of 0.8. Both flights saw a rapid decrease in performance moving from a probability threshold of 0.8 to 0.9.

Figure 10 shows the visualizations for both the 200 ft and 300 ft flights at probability thresholds of 0.1 and 0.8. The largest deviation is seen in the 200 ft flight along a turn, which was then able to recover and continue on the appropriate path.



Figure 8. Visualization of the mapping results using Flight Path #1 with two different IoU Association Thresholds 0.1 on top and 0.8 on bottom. Note that the bounding boxes are visualized on the available satellite imagery taken on a different day than the flights. The bounding boxes show the car locations seen during the flight, not those seen on the satellite imagery.

4.4 Comparison to Keypoint Localization

To better assess the key-region performance, this study directly compares the proposed algorithm to a re-implementation of the state-of-the-art keypoint based localization (Wei and Yilmaz, 2023). For this study, the two algorithms were run on Flight Path #2 flown at 300 ft with varying levels of Gaussian blur to simulate image degradation. Table 3 shows the results for this study. Method refers to either the keypoint (KP) or key-region (KR) based localization. The reference map size (Ref. Size) refers to the number of landmarks used for that flight's localization. The keypoint based localization used two difference sized reference maps controlled by the non-maximum suppress-

Probability Thresh.	Flight Height (ft)	Mean Error (m)
0.1	200	3.51
0.2	200	3.52
0.3	200	3.36
0.4	200	3.36
0.5	200	4.97
0.6	200	4.97
0.7	200	6.91
0.8	200	6.91
0.9	200	53.08
0.1	300	2.26
0.2	300	2.26
0.3	300	2.26
0.4	300	2.26
0.5	300	2.24
0.6	300	2.24
0.7	300	2.23
0.8	300	2.23
0.9	300	18.67

Table 2. IoU Association Threshold positioning results.



Figure 9. Visualization of the localization results for the 200 ft and 300 ft flights over Flight Path #2.

sion (NMS) of 1 and 25, where they produced 6115 and 1224 keypoints respectively. The key-region landmarks are the extracted object identities, which for this map was 37. The blur refers to the neighborhood of influence for Gaussian blur.

The key-region based localization proved significantly more resilient to image degradation compared to the key-point based localization, maintaining a consistent mean positioning error of 2.23 m. The proposed algorithm also outperformed the compared algorithm when it was limited to the smaller reference map size. When given the higher density of keypoints, the larger reference map, the key-point based localization was more accurate with a mean error of 1.83 m. The processing time of the key-region based localization was slower at 1.17 FPS compared to 3.05 FPS and 1.90 FPS of the key-point based method.

Method	Ref. Size	Blur	Mean Error	FPS
KR	38	0	2.24	1.17
KR	38	3	2.24	1.17
KR	38	5	2.24	1.17
KR	38	7	2.24	1.16
KP	1224	0	2.23	3.00
KP	1224	3	2.25	3.08
KP	1224	5	2.31	3.06
KP	1224	7	2.31	3.07
KP	6115	0	1.83	1.85
KP	6115	3	1.86	1.92
KP	6115	5	1.89	1.89
KP	6115	7	1.92	1.92

Table 3. Results for comparison of key-region (KR) versus keypoint (KP) localization methods.

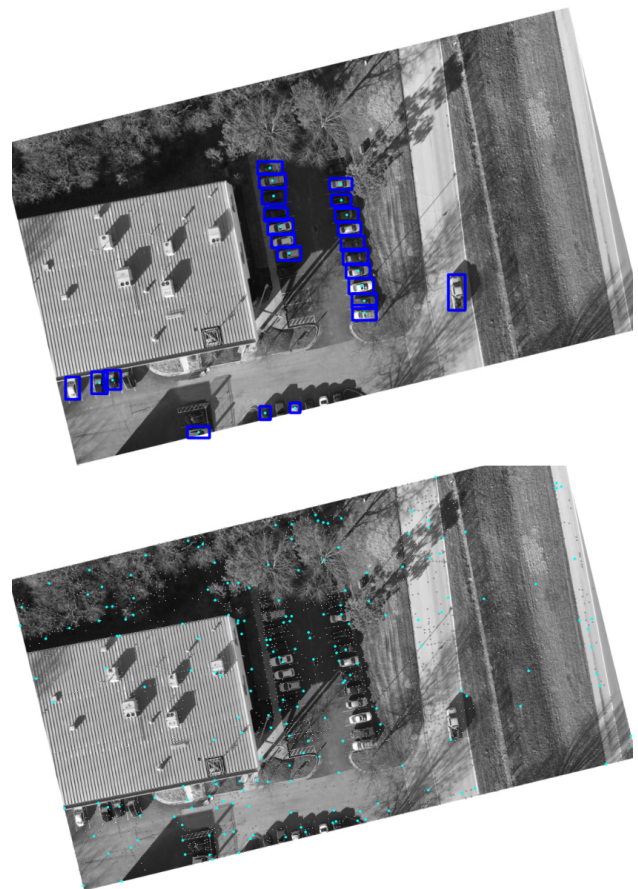


Figure 10. Visualization of the key-region and keypoint based localization flights over Flight Path #2 with a Gaussian blur of 3.

5. Conclusion

Key-region based visual navigation shows strong feasibility from the limited studies completed in this work. With the proper parameters, the system achieved a mean localization error of 2.23 m and 3.36 m. This performance is comparable to current GPS based and keypoint based localization methods producing localization accuracies between 0.50 m and 3.4 m. The key-region based approach proved more resilient to image degradation simulated through increasing Gaussian blur. Overall, key-region based visual navigation showed comparable localization performances to similar alternatives with improved resiliency and higher understandability.

References

- Akyon, F. C., Cengiz, C., Altinuc, S. O., Cavusoglu, D., Sahin, K., Eryuksel, O., 2021. SAHI: A lightweight vision library for performing large scale object detection and instance segmentation.
- Chen, Y.-E., Liew, H.-H., Chao, J.-C., Wu, R.-B., 2023. Decimeter-Accuracy Positioning for Drones Using Two-Stage Trilateration in a GPS-Denied Environment. *IEEE Internet of Things Journal*, 10(9), 8319-8326.
- Czyża, S., Szuniewicz, K., Kowalczyk, K., Dumalski, A., Ogrodniczak, M., Zieleniewicz, , 2023. Assessment of Accuracy in Unmanned Aerial Vehicle (UAV) Pose Estimation with the REAL-Time Kinematic (RTK) Method on the Example of DJI Matrice 300 RTK. *Sensors*, 23, 2092.
- Elkhrachy, I., 2021. Accuracy Assessment of Low-Cost Unmanned Aerial Vehicle (UAV) Photogrammetry. *Alexandria Engineering Journal*, 60(6), 5579-5590. <https://www.sciencedirect.com/science/article/pii/S1110016821002544>.
- Frosi, M., Gobbi, V., Matteucci, M., 2023. OSM-SLAM: Aiding SLAM with OpenStreetMaps priors. *Frontiers in Robotics and AI*, 10, 1064934.
- Gupta, A., Fernando, X., 2022. Simultaneous Localization and Mapping (SLAM) and Data Fusion in Unmanned Aerial Vehicles: Recent Advances and Challenges. *Drones*, 6(4). <https://www.mdpi.com/2504-446X/6/4/85>.
- Jametoni, F., Saputra, D. E., 2021. A study on autonomous drone positioning method. *2021 Sixth International Conference on Informatics and Computing (ICIC)*, 1–5.
- Karnes, M., Yilmaz, A., 2023. Omni-Modeler: Rapid Adaptive Visual Recognition with dynamic learning. *Signal and Image Processing: An International Journal*, 14(4/5), 01–12.
- Li, Z., Zhao, C., Wang, J., Hou, X., Hu, J., Pan, Q., Jia, C., 2022. High-accuracy robust slam and real-time autonomous navigation of uav in gnss-denied environments. L. Yan, H. Duan, X. Yu (eds), *Advances in Guidance, Navigation and Control*, Springer Singapore, Singapore, 1099–1108.
- Mazurek, P., Hachaj, T., 2021. SLAM-OR: Simultaneous Localization, Mapping and Object Recognition Using Video Sensors Data in Open Environments from the Sparse Points Cloud. *Sensors*, 21(14). <https://www.mdpi.com/1424-8220/21/14/4734>.
- Mráz, E., Rodina, J., Babinec, A., 2020. Using fiducial markers to improve localization of a drone. *2020 23rd International Symposium on Measurement and Control in Robotics (ISMCR)*, 1–5.
- Patrik, A., Utama, G., Gunawan, A., Chowanda, A., Suroso, J., Shofiyanti, R., Budiharto, W., 2019. GNSS-based navigation systems of autonomous drone for delivering items. *Journal of Big Data*, 6.
- Pi, Y., Nath, N., Bezhadan, A., 2020. Deep neural networks for drone view localization and mapping in gps-denied environments. 1–16.
- Supriyono, H., Akhara, A., 2021. Design, building and performance testing of GPS and computer vision combination for increasing landing precision of quad-copter drone. *Journal of Physics: Conference Series*, 1858, 012074.
- Wang, K., Kooistra, L., Pan, R., Wang, W., Valente, J., 2024. UAV-based simultaneous localization and mapping in outdoor environments: A systematic scoping review. *Journal of Field Robotics*.
- Wei, J., Yilmaz, A., 2023. A Visual Odometry Pipeline for Real-Time UAS Geopositioning. *Drones*, 7(9). <https://www.mdpi.com/2504-446X/7/9/569>.
- Zeng, W., ying Xiao, Z., 2024. Few-shot learning based on deep learning: A survey. *Mathematical Biosciences and Engineering*, 21(1), 679-711. <https://www.aimspress.com/article/doi/10.3934/mbe.2024029>.