

# Multiview Aerial Visual Recognition (MAVREC): Can Multi-view Improve Aerial Visual Perception?

Anonymous CVPR submission

Paper ID 16869

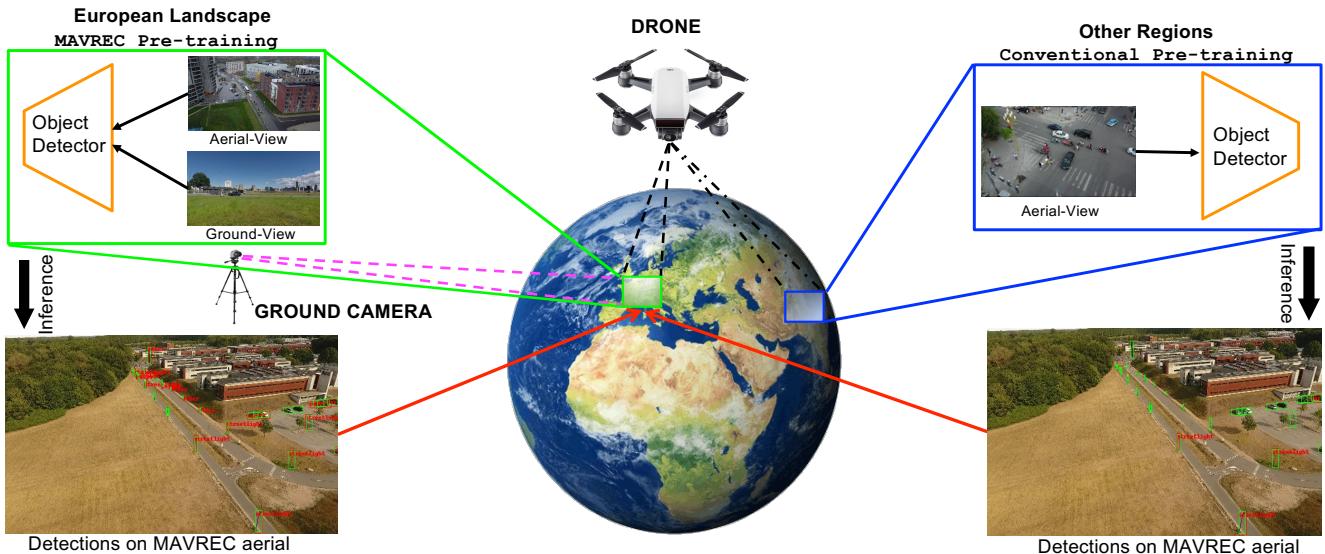


Figure 1. Illustration of the geography-aware model using our proposed MAVREC dataset (green box) collected in the rural and urban European landscape vs. the conventional aerial object detector (blue box) pretrained only on aerial images from VisDrone [88] captured in Asia. The conventional approach fails to precisely detect aerial objects from the MAVREC dataset. In contrast, our object detector pretrained on the ground and aerial images from the MAVREC dataset contextualize the object proposals of that specific geography and enhance the aerial visual perception, thus outperforming other object detectors pre-trained on popular ground-view dataset (MS-COCO [44]) or other aerial datasets collected from different geographies; also, see Figure 6.

## Abstract

Despite the commercial abundance of UAVs, aerial data acquisition remains challenging, and the existing Asia and North America-centric open-source UAV datasets are small-scale or low-resolution and lack diversity in scene contextuality. Additionally, the color content of the scenes, solar-zenith angle, and population density of different geographies influence the data diversity. These two factors conjointly render suboptimal aerial-visual perception of the deep neural network (DNN) models trained primarily on the ground-view data, including the open-world foundational models.

To pave the way for a transformative era of aerial detection, we present **Multiview Aerial Visual RECognition** or **MAVREC**, a video dataset where we record synchronized scenes from different perspectives — ground camera

and drone-mounted camera. MAVREC consists of around 2.5 hours of industry-standard 2.7K resolution video sequences, more than 0.5 million frames, and 1.1 million annotated bounding boxes. This makes MAVREC the largest ground and aerial-view dataset, and the fourth largest among all drone-based datasets across all modalities and tasks. Through our extensive benchmarking on MAVREC, we recognize that augmenting object detectors with ground-view images from the corresponding geographical location is a superior pre-training strategy for aerial detection. Building on this strategy, we benchmark MAVREC with a curriculum-based semi-supervised object detection approach that leverages labeled (ground and aerial) and unlabeled (only aerial) images to enhance the aerial detection. We publicly release the MAVREC dataset: <https://mavrec.github.io>.

015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029

## 030 1. Introduction

031 Object detection and tracking employing UAV (or drone)-  
032 based aerial videos are essential in many downstream applica-  
033 tions, such as autonomous driving [15, 52], robotics [78],  
034 environmental monitoring [59], infrastructure inspection  
035 [13], developing livable and safe communities [6, 30, 86],  
036 a few to name. Despite having many crucial applications,  
037 most visual perception models focus on ground-view im-  
038 ages. This bias results in suboptimal performance when  
039 these models are applied to an aerial perspective — a dis-  
040 crepancy due to a domain shift precipitated by the viewpoint  
041 transfer. We hypothesize that the basis of this disparity is  
042 primarily twofold:

043 First, the lack of diversity in the current aerial datasets.  
044 Modern DNN-based visual models are data-hungry. How-  
045 ever, aerial data collection is intricate due to UAV flight  
046 regulations and safety protocol, atmospheric turbulence, and  
047 many more [31]. The existing open-source UAV datasets  
048 [16, 17, 21, 49, 51, 54, 67, 88] are either small-scale, or  
049 low-resolution, and collected primarily in the urban pasture  
050 across Asian and North American geographies. These fac-  
051 tors contribute to inadequacy in diverse dataset properties  
052 and hinder training large DNNs for aerial visual perception.

053 Second, substandard generalizability of the existing aerial  
054 visual models across different geographic locations. An ob-  
055 ject detection model trained on datasets captured from South  
056 Asia underperforms when deployed on videos captured with  
057 European landscapes, characterized by semi-rural pastures  
058 and lots of greenery; see Figure 1. Research shows that dif-  
059 ferent geographic factors, including latitude influences the  
060 population density [8, 37, 38], hence, the *color-content of the*  
061 *scenes*<sup>1</sup>, their complexities, and density and interactions of  
062 the foreground objects. These seemingly low-key factors di-  
063 rectly affect the data captured from a drone-mounted camera,  
064 and the previous studies never considered the inter-domain  
065 inference quality of the DNN models trained on these data.

066 To pave the way for a transformative era of aerial visual  
067 perception models, we introduce **Multiview Aerial Visual**  
068 **R****E****C****ognition** dataset, **MAVREC**, which uniquely captures  
069 time-synchronized aerial and ground view data. MAVREC  
070 is collected with consumer-grade handheld cameras (smart-  
071 phones and GoPro) and drone-mounted cameras, consists of  
072 around 2.5 hours of industry-standard 2.7K resolution video  
073 sequences, more than 0.5 million frames, covering rural and  
074 urban pastures during spring and summer in high-latitude  
075 European geographies. It makes *MAVREC the largest ground*  
076 *and aerial-view dataset, and the fourth largest among all*  
077 *drone-based datasets (edited and unedited) across all modal-*  
078 *ties and tasks that ever existed*; see Table 1.

079 In this paper, we rigorously assess our hypothesis and  
080 explore interesting properties of object detection in aerial

<sup>1</sup>European vehicles are comprising of mainly three colors [5]; also, see B.2 for an analysis.

081 images while evaluating MAVREC in a supervised setting.  
082 We find that aerial object detection considers contextual in-  
083 formation of the landscape and thus is different from general  
084 object detection in ground view. An object detection model  
085 trained on a specific geographic location may not generalize  
086 to other geographic locations. Thus, aerial object detection  
087 models require learning geography-aware representation of  
088 aerial images. Finally, building on this hypothesis, we bench-  
089 mark MAVREC with a curriculum-based semi-supervised  
090 object detection approach that leverages labeled and unla-  
091 beled images to enhance the detection performance from an  
092 aerial perspective.

093 We summarize our key technical contributions as follows:

- We introduce MAVREC, which to date represents the most extensive dataset integrating time-synchronized ground and aerial images captured in the European landscape; §3.
- Through benchmarking MAVREC in supervised and semi-supervised settings, we expose the proclivity of existing pre-trained object detectors to exhibit bias toward data captured from ground perspectives; §4.
- We propose a curriculum-based semi-supervised object detection method. Its superior performance shows the importance of training these types of models with ground-view images to learn geography-aware representation; §4.2.

## 2. Related work

105 **UAV-based datasets.** The last decade witnessed a surge in  
106 UAV-based video and image datasets. We list some open-  
107 source UAV datasets, curated since 2016, and group their  
108 key features according to their downstream tasks.

109 *VisDrone* [88] is the most widely used drone dataset for  
110 aerial image object detection. It is recorded from 14 cities  
111 in China with various drone-mounted cameras, consists of  
112 10 object categories, and segregated into four task-specific  
113 sub-datasets: (a) Image Object Detection (10,209 images),  
114 (b) Video Object Detection (96 videos, 40,001 images), (c)  
115 Single-Object Tracking (139,276 images), and (d) Multi-  
116 Object Tracking (40,000 images). *Campus* [63], is the largest  
117 aerial dataset for multi-target tracking, activity comprehen-  
118 sion, and trajectory prediction, focuses solely on the uni-  
119 versity campus, in contrast to our MAVREC. *UAVDT* [21]  
120 dataset consists of 80,000 frames and 3 subsets, focusing  
121 on single and multi-object detection and tracking, under dif-  
122 ferent weather condition, lighting, and altitude of the drone.  
123 *MOR-UAV* [51] is an aerial dataset consisting of 10,948  
124 images, all annotated, designed for moving object detec-  
125 tion under various challenges, such as illumination, camera  
126 movement, etc. *UAV123* [54] is a low-altitude aerial dataset  
127 consisting of 112,578 fully-annotated images across 123  
128 video sequences (simulated and recorded), designed for ob-  
129 ject tracking, with a subset intended for long-term aerial  
130 tracking. *MDOT* [87] is a *multi-drone based single object*  
131 *tracking dataset* with 259,793 frames across 155 groups of

Table 1. State-of-the-art UAV-based datasets since 2016 in chronological order. For viewpoints, G denotes *ground-view*, A denotes *aerial-view*, and AG denotes both. Thermal IR datasets are not included.

Dataset	Total Frames	Resolution	Total Annotations	Instances per Annotated Frame	Categories	Viewpoints	Region	Year
Campus [63]	929,499	1400 × 2019	19,564	0.02	6	Single (A)	North America	2016
UAV123 [54]	110,000	720 × 720	110,000	1.0	6	Multi (A)	Middle East	2016
CarFusion[61]	53,000	1, 280 × 720	—	—	4	Multi	North America	2018
DAC-SDC [82]	150,000	640 × 360	NA	NA	12	Single	Asia	2018
UAVDT [21]	80,000	1080 × 540	841,500	10.52	3	Single	Asia	2018
MDOT [87]	259,793	—	—	—	9	Multi (A)	Asia	2019
Visdrone DET [88]	10,209	3840 × 2160	471,266	53.09	10	Single (A)	Asia	2019
Visdrone MOT [88]	40,000	3840 × 2160	1,527,557	45.83	10	Single (A)	Asia	2019
DOTA[79]	2,806	4000 × 4000	188,282	67.09	15	Single (A)	Multiple	2019
DOTA V2.0 [20]	11,268	4000 × 4000	1,793,658	159.18	18	Single (A)	Multiple	2021
MOR-UAV [51]	10,948	1280 × 720, 1920 × 1080	89,783	8.20	2	Single	Asia	2020
AU-AIR [17]	32,823	1920 × 1080	132,034	4.02	8	Multi	Europe	2020
UAVid [49]	300	3840 × 2160	—	—	8	Single	Europe	2020
MOHR [84]	10,631	7360 × 4192, 8688 × 5792	90,014	8.47	5	Multi (A)	Asia	2021
<b>MAVREC (This paper)</b>	<b>537,030</b>	<b>2700 × 1520</b>	<b>1,102,604</b>	<b>50.01</b>	<b>10</b>	<b>Multi (AG)</b>	<b>Europe</b>	<b>2023</b>

video clips, and 10 different annotated attributes. *Au-Air* [17] is a medium scale, multi-sensor, aerial data designed for real-time object detection, with the aim of bridging the gap between computer vision and robotics. *DAC-SDC* [82] is a single-object detection dataset with 150,000 images collected from *DJI* [7] with 12 categories. *DOTA* [79] is an aerial dataset (2,806 images, 15 categories) for object detection in earth vision. *DOTA V2.0* [20], an upgraded version of *DOTA*, is a single-view dataset collected from Google Earth, GF-2 satellite and aerial imagery (11,268 images, 18 categories) for object detection. The *UG<sup>2</sup>+ Challenges* provide *A2I2-Haze* [56], the first real haze dataset, consisting of 229 pairs of hazy and clean images (197 training pairs, 32 testing pairs) from 12 videos, focusing on detection in visually degraded environments with smoke and haze with mutually exclusive aerial and ground images.

In an orthogonal line of work, *GRACO* [91] is a multimodal dataset for synchronized ground and aerial collaborative simultaneous localization and mapping (SLAM) algorithms (6 ground and 8 aerial sequences collected in China within a university campus) by a group of ground and aerial robots equipped with light detection and ranging (LiDAR), cameras, and global navigation satellite/inertial navigation systems (GNSS/INS) that capture images at 20Hz with 2182.52 and 2675.54 seconds duration in the ground and aerial, respectively. *S3E* [23] is a multimodal dataset for collaborative SLAM, consists of 7 outdoor and 5 indoor synchronized ground sequences, each longer than 200 seconds, and collected in 5 locations within a university campus in China. *DVCD18K* [12] is a cinematographic dataset (with corresponding camera paths) consisting of 18,551 edited drone clips spanning 44.3 hours.

Our proposed *MAVREC* is inherently different from the above datasets, because: (a) compared to other small-scale (e.g., UAVid, DOTA, MOR-UAV, AU-AIR), and low-resolution datasets (e.g., UAV123, UAVDT, DroneVehicle, BIRDSE), *MAVREC* is the first-ever *large-scale*, un-

scripted, multi-viewpoint video dataset (fourth largest among all UAV-based datasets ever after DVCD18K, Highway-drone, and Campus; Campus lacks object boundaries) recorded in *industry-standard* 2.7K resolution; (b) its multi-viewpoint presents the same scenes through the lens of one or more ground cameras, and a medium altitude (flight height 25–45 meters, compared to low or high-altitude datasets, e.g., UAV123 with flight height 5–25 meters, MOHR with flight height 200 meters or above) drone-mounted camera (this perspective is unique compared to the existing multiview drone datasets, e.g., MDOT, UAV123, MVDTD, MCL) to have balanced variations of small and medium objects from both perspectives; (c) its high variance in object distribution across different scenes is complementary to datasets like VisDrone where object detection is relatively straightforward due to their biased object distribution (dense), reflecting its demographic characteristics; (d) *MAVREC* is the first multi-view, drone-based dataset curated in the wild, with a central focus on object detection. *GRACO* and *S3E* are multi-view but confined to university campuses and used for SLAM algorithms; *S3E* does not incorporate drone-based data acquisition. Alongside, *A2I2-Haze* dataset a tiny subgroup of [56] challenges, consists of mutually exclusive, non-synchronized aerial and ground images, while *DVCD18K* is an *human-edited cinematographic dataset with drone camera paths*, and vastly different from *MAVREC* in many aspects. §3.2 explains more unique challenges of *MAVREC*. There are UAV-based datasets with downstream tasks primarily orthogonal to *MAVREC*. For completeness, we list some UAV-based datasets for action detection, counting, geo-localization, 3D reconstruction, and benchmarking in §A; also, see [78].

(iii) **Object detection.** CNN-based object detectors are divided into two categories: two-stage and one-stage. Two-stage detectors such as RCNN [26], Fast RCNN [25], Faster RCNN [62], employ a class-agnostic region proposal module followed by simultaneously regressing the object bound-



Figure 2. Different sample scenes (with annotation) from our dataset; the first row is the aerial-view, second row presents the same scenes from a ground camera. See more sample frames in §B, Figure 9.

aries and their classes. In contrast, one-stage detectors like SSD [47], YoloV4 [72], YoloV6 [40], YoloV7 [73], YoloX [24], FCOS [68], directly predicts the image pixels as objects, leading to models that offer fast inference. Recently, by using neural architecture search, Yolo-NAS [11] outperforms previous Yolo models in real-time object detection. With the success of transformers, DETR [18] was the first transformer-based, end-to-end object detector. Following this, Deformable-DETR (D-DETR) [89] introduces a sparse attention module, computationally  $6\times$  faster than DETR, and robust in detecting small objects. The majority of object detectors designed for aerial imagery draw upon the foundational principles established by these aforementioned popular object detectors [75, 80, 81]. Along this line, TPH-YoloV5 [90] combines YoloV5 with a transformer prediction head to solve the varying object scales and motion blur for drone-captured scenarios. As a result, our analysis utilizes the MAVREC dataset to benchmark these well-established methods, prioritizing factors such as fast inference, high precision, and the effective detection of small-scale objects.

### 226 3. MAVREC Dataset

227 In this section, we start with the data acquisition process;  
228 and then explain annotation, statistical attributes, and unique  
229 challenges of MAVREC.

#### 230 3.1. General setup

231 **Recording set-up.** We conduct the recording in public  
232 spaces in compliance with the European Union’s drone  
233 safety and Scandinavian video surveillance regulations; see  
234 §D for detailed discussion on reproducibility, licensing,  
235 privacy, safety, maintenance plan and broader impact of the  
236 dataset. We record our dual-view aerial-ground dataset with  
237 a drone-mounted camera (DJI Phantom 4, DJI mini 2) and  
238 a consumer-grade static ground camera (GoPro Hero 4, Go-  
239 Pro Hero 6, iPhone 11 and 13-Pro) placed on a tripod; see  
240 details in Table 5. The drone is kept semi-static, hovering ap-  
241 proximately 25–45 meters above the ground; see the relative  
242 positions and viewing angles of the drone and the ground  
243 camera in Figure 3. Based on that, we identify three record-

ing scenarios (P1, P2, and P3). In P3, we better capture the  
244 objects as the drone gets a wider viewing angle. However,  
245 we keep all views *not to amplify biases* from any particular  
246 view. For some recordings in the city center, railroad, or  
247 crowded intersections, we were unable to operate a drone  
248 due to the UAV-flight regulations; hence, we used a user-  
249 grade handheld camera set-up in the balcony of a high-riser  
250 to capture aerial views.

**Recording locations and scenes.** To avoid locational bias,  
252 we collected our data in 11 different geographical locations  
253 (European outdoors, rural and urban) with mixed pastures,  
254 in spring and summer (with the sun hitting the cameras from  
255 different angels), and when there is an encyclopedic spec-  
256 trum of green and yellow intertwined in the background; see  
257 Figures 3 and 2 (also, see B.2 for an analysis). We choose the  
258 parking lots, and busy traffic intersections in the city, during  
259 the peak traffic hours to create more nuanced and complex  
260 interactions, in which multiple foreground objects are inter-  
261 acting and creating enormous visual challenges. Alongside,  
262 we choose harbor, single-lane roads in the countryside, as-  
263phalt roads, and bi-cycle lanes, in moderate traffic conditions,  
264 to collect simple scenarios which might have sparse to dense  
265 foreground objects (see sample frames in Figure 2).

**Alignment of dual-views.** Human operators simultaneously  
267 record the scenes from dual views; although a minute time-  
268 lapse is unavoidable. Consequently, after recording, clips  
269 are loaded into the QuickTime player, and a human opera-  
270 tor manually synchronizes the frames to alleviate the time  
271 lapses. We note that for 12.5% of the clips (22.3% of im-  
272 age frames), we captured an extra ground view. Thus, these  
273 video clips offer three perspectives in total. The inclusion of  
274 a second ground camera aims to enhance ground-to-ground  
275 representation. Although these videos have not been utilized  
276 in this paper, they are preserved for future work.

**Annotation and categories.** MAVREC, as highlighted previ-  
278 ously, stands as one of the largest drone datasets, encompass-  
279 ing millions of objects within its distribution. However, anno-  
280 tating each object is a resource-intensive task. Inspired by the  
281 recent success of the semi-supervised learning paradigm in  
282

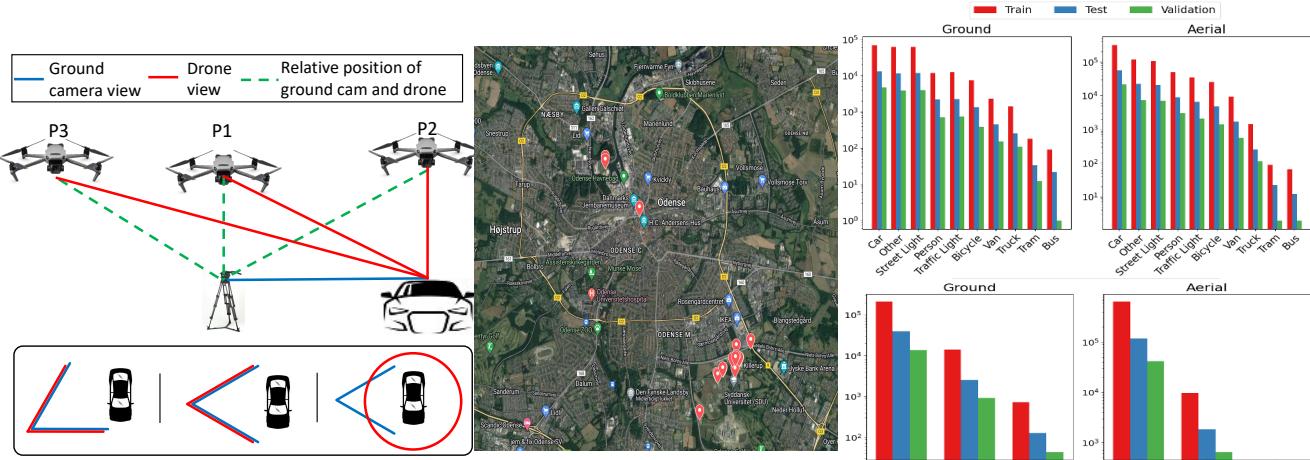


Figure 3. **Left:** Recording instances are divided into three different scenarios (P1, P2, and P3) based on the relative positions and the field-of-view (FOV) of the ground camera and the drone. The drone operates directly on top of the object (P2), and two oblique views—directly on top of the ground camera (P1), and behind the ground camera (P3). **Middle:** Recording locations as red dropped pins on the google map’s sattelite view. **Right (top):** Total numbers of objects in each category in the ground and aerial view. **Right (bottom):** Number of large, medium, and small objects in the test and validation sets; aerial view has no *large* object annotation.

283 the computer vision community [19, 39, 48, 53, 74, 76], we  
 284 split the videos from different views into two categories; an  
 285 annotated set and an unannotated set. After pre-processing,  
 286 we select the first 30 seconds of the synchronized videos  
 287 and annotate the frames through a semi-automatic, open-  
 288 source annotation platform by Intel, called CVAT [2], and  
 289 leave the rest of the video unannotated; see CVAT inter-  
 290 face in Figure 10. We provide an annotation interface  
 291 with 10 categories in CVAT: tram, truck, bus, van,  
 292 car, bicycle, person, street light, traffic  
 293 light, and other. In other category, we annotated  
 294 objects that share visual similarities with objects from the  
 295 remaining categories; e.g., blocks of concrete from aerial view  
 296 might look like cars, or white divider and marker posts from  
 297 aerial view might look like a person with a white T-shirt, and  
 298 so on. We created this category for the models to learn to  
 299 disambiguate the *look-alike* objects from different categories.  
 300 The in-built tracker in CVAT tracks an object through mul-  
 301 tiple frames. We annotated by skipping forward 10 frames;  
 302 thus speeding up the annotation process. Nevertheless, to  
 303 ensure high annotation quality, a human annotator reviews  
 304 each frame and 5 human non-annotator have reviewed the  
 305 dataset annotations. The annotation reviewers check that the  
 306 bounding box encapsulates the object or its parts, has mini-  
 307 mal overlap with other objects, and that all instances of the  
 308 class in the frame are labelled. This process is performed on  
 309 600 annotated images randomly samples from the annotated  
 310 data. During the review process, we find that the error is  
 311 around 6% which is comparable to the benchmark datasets  
 312 in this domain. Similar to other benchmarks [44, 88], anno-  
 313 tated frames are assembled into COCO-.json format to give  
 314 a unique identifier for each object class.

### 3.2. Structuring, statistics, and challenges

This section discusses the size, statistical properties, and challenges of the training distribution of MAVREC. We focus on two key points: (i) distribution of different categories, and (ii) distribution of the annotated object size.

**Structuring the dataset.** We divided the annotated data from both views into three subsets—train, validation, and test sets. To ensure the distributions of the different objects are approximately the same throughout these three sets, we split each video sequence into three fragments, and then randomly select samples for each set.

**Distribution of different categories.** We show the distribution of categories from both views in Figure 3; also, see Figure 11. MAVREC contains over 1.1 million bounding box annotations in both views combined, rendering  $\sim 50.01$  annotations per frame; see Figure 3 and details in Table 6. The distribution is *long-tailed* where cars are more frequent than trams and buses. The slight inconsistency in the object distribution from both views is natural as some recordings were conducted with the P3 setup, and in this setup, the drone has a wider viewing angle than the ground camera.

**Object size distribution.** To better illustrate the challenges in MAVREC, we divide the object sizes present in the videos into *three* categories: small ( $< 32 \times 32$  pixels), medium (lies inclusively between  $32 \times 32$  and  $96 \times 96$  pixels), and large ( $> 96 \times 96$  pixels). Figure 3 (also, see Figure 12) presents the number of annotated object sizes in both views. Large objects, such as trams, buses, and trucks, are present in fewer frames compared to the other objects. Also, the drone is maneuvered at a higher altitude, and the aerial view has a higher percentage of small objects compared to the ground

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346 view, creating a natural bias in object sizes. We also observe  
347 that the distribution for the split into the train, validation,  
348 and test set has almost an equal distribution of the different  
349 object sizes for both views; see Figure 12-(c) for distribution  
350 for the object sizes.

351 **Unique properties of MAVREC.** MAVREC contains typi-  
352 cal outdoor activities characterized by real-world properties  
353 like long-tail distribution, objects with similar appearance,  
354 viewpoint changes, varying illumination, etc. Additionally,  
355 MAVREC exhibits some unique properties, not found in  
356 other datasets: (i) Ground view contains occluded objects.  
357 Nevertheless, these objects can be recovered due to the wide  
358 aerial FOV. This *dual-view feature of the MAVREC* has the  
359 potential to offer a wide range of solutions for scenes with  
360 occlusion, which remains a significant challenge in video  
361 surveillance. (ii) *MAVREC’s color distribution* reflects Euro-  
362 pean demographics, which may influence object detection  
363 algorithms that incorporate scene-contextual information,  
364 particularly those pre-trained on general object detection  
365 datasets; see a comparison in Figure 13. (iii) Historically,  
366 vehicle color distributions vary across Europe, North Amer-  
367 ica, and the Asia-Pacific; see Figure 8. The existing datasets  
368 collected in Asia and North America appear to be more col-  
369 orful. E.g., in 2021, Europe’s top car colors were gray (27%),  
370 white (23%), and black (22%), contrasting with North Amer-  
371 ica’s gray (21%), black (20%), blue (10%), red (10-11%),  
372 and silver (10%), and China’s predominance of white (50%)  
373 and brown (10%) cars [4]. (iv) MAVREC was collected  
374 at *high latitudes*. The elevation of the sun in these areas  
375 (see Figure 7) during the peak traffic times is high, creat-  
376 ing a *mirage-like* reflection on one of the sensors in many  
377 scenes, thereby, causing significant disparities between the  
378 two views. The second column of rows 1 and 2 in Figure  
379 2 shows this effect. (v) The aerial perspective inherent in  
380 MAVREC leads to *small objects* inclusion; their presence is  
381 susceptible to miss-detection by detection algorithms. (vi)  
382 MAVREC is characterized by both *sparse* and *dense*  
383 distribution of objects. Our empirical findings suggest that such a  
384 large disparity in object distribution presents challenges in  
385 training object detectors, compared to scenes exhibiting only  
386 uniformly dense annotations.

## 387 4. Baselines and evaluation

388 This section presents the benchmarking results on MAVREC  
389 in supervised and semi-supervised settings. We also present  
390 our observations concerning the prevailing trends in object  
391 detectors employed on aerial images.

392 **Datasets and evaluation metric.** For supervised and semi-  
393 supervised benchmarking with MAVREC, we use a total of  
394 8,605 labeled frames, and at most 8,605 unlabeled frames  
395 from each view at training. The validation and test set (ex-  
396 tracted from disjoint video sequences) for each view contain  
397 805 and 1,614 annotated images, respectively. We evaluate

398 the models with the widely used metric for object detection,  
399 mean average precision (mAP) [44]; see a detailed discus-  
400 sion in §C.2.

401 **Object detector baselines.** For *supervised benchmarking*,  
402 we use CNN-based YoloV7 [73], and transformer-based  
403 DETR [18] and D-DETR [89]. Additionally, we use Yolo-  
404 NAS [11]. For *semi-supervised benchmarking*, we propose a  
405 curriculum based semi-supervised baseline using D-DETR.  
406 We provide the implementation details and computing en-  
407 vironment in the §C.1; we refer to Tables 7 and 8 for other  
408 model specific implementation details.

### 409 4.1. Supervised benchmarking

410 Table 2 presents the supervised baselines results on  
411 MAVREC dataset for both ground and aerial perspectives.  
412 Despite an equal number of training samples from differ-  
413 ent views, we observe that all the baselines exhibit su-  
414 perior performance on the ground perspective compared to  
415 the aerial perspective. This discrepancy highlights the chal-  
416 lenge associated with object detection in aerial views due to  
417 their smaller sizes, as indicated by the AP<sub>S</sub> metric. Notably,  
418 YoloV7 demonstrates the best performance on aerial images,  
419 while D-DETR pre-trained on MSCOCO surpasses other  
420 models on the ground view. Interestingly, Yolo-NAS, which  
421 surpasses other Yolo-based detectors on ground images ac-  
422 cording to [11], exhibits lower performance than YoloV7  
423 on aerial images indicating that the learned Yolo-NAS ar-  
424 chitecture is suboptimal for aerial images. We show some  
425 qualitative results in Figure 16.

426 **Can ground-view images improve object detection in  
427 aerial perspective?** To answer this, we trained D-DETR  
428 and YoloV7 by augmenting the existing aerial-view sample  
429 set with ground-view samples. We achieve this by simply  
430 adding the two sets of aerial- and ground-view samples along  
431 with their corresponding annotations. Our findings demon-  
432 strate that the inclusion of ground-view samples substantially  
433 improves the object detection performance.

434 Graphs presented in figure 4 illustrates that D-DETR  
435 outperforms the CNN-based YoloV7 when the extra ground-  
436 view samples enrich the training distribution. While YoloV7  
437 requires an equal number of ground-view samples as aerial  
438 samples to achieve its peak performance, D-DETR achieves  
439 a relative improvement upto 270% with a subset of ground-  
440 view samples (~ 2K ground-view images). Interestingly, fur-  
441 ther augmentation of ground-view images during D-DETR  
442 training does not enhance its performance, indicating the  
443 sensitivity of D-DETR’s training process to ground-view  
444 image sampling.

445 Thus, experimental analysis of MAVREC within a super-  
446 vised framework suggest that (i) CNN-based models, such as  
447 Yolo, demonstrated superior performance over transformer-  
448 based models, like D-DETR, when trained from scratch.  
449 However, transformer based model outperforms when pre-  
450 trained on large-scale ground images in MSCOCO. This

Table 2. Supervised benchmark of MAVREC. D-DETR\* denotes a MSCOCO pre-trained D-DETR.

Trained	Validation Set								Test Set											
	DNN				Ground				Aerial				Ground				Aerial			
	Models		AP	AP <sub>50</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP	AP <sub>50</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP	AP <sub>50</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP	AP <sub>50</sub>	AP <sub>S</sub>	AP <sub>M</sub>		
DETR	21.8	36.9	21.9	23.9	24.9	39.7	27.6	45.3	20.8	35.4	21.3	24.0	23.6	40.1	23.4	44.9				
D-DETR	27.5	51.4	28.1	43.7	13.1	28.3	14.2	38.1	18.2	46.8	17.9	36.0	10.3	25.0	10.1	29.4				
D-DETR*	<b>59.6</b>	<b>82.7</b>	<b>59.7</b>	<b>79.6</b>	31.0	<b>61.7</b>	31.7	55.1	<b>58.6</b>	<b>81.4</b>	<b>59.0</b>	<b>80.2</b>	<b>33.2</b>	<b>61.9</b>	<b>31.5</b>	<b>51.0</b>				
Yolo-NAS (L)	41.4	61.7	36.8	72.9	30.3	49.8	29.2	61.5	41.2	63.4	37.8	74.3	27.0	43.3	25.9	58.0				
YoloV7	45.6	72.1	40.6	74.9	<b>31.3</b>	57.7	<b>34.2</b>	<b>61.2</b>	45.0	72.5	42.4	74.4	31.9	58.8	31.4	<b>63.1</b>				

Training Protocol	Pre-training	Test Set			
		AP	AP <sub>50</sub>	AP <sub>S</sub>	AP <sub>M</sub>
Trained from scratch	X	10.3	25.0	10.1	29.4
Grounding-DINO	O365, GoldG, Cap4M	20.4	40.9	18.6	32.5
FT on MAVREC Aerial view	Visdrone [88]	20.9	41.9	20.6	43.8
FT on MAVREC Aerial view	MS-COCO [44]	33.2	61.9	31.5	51.0
FT on MAVREC Aerial view	MAVREC Ground view	<b>44.8</b>	<b>71.5</b>	<b>42.9</b>	<b>72.4</b>

Table 3. Object detection using D-DETR [89] on the aerial view images of the proposed MAVREC dataset; FT indicates finetuning.

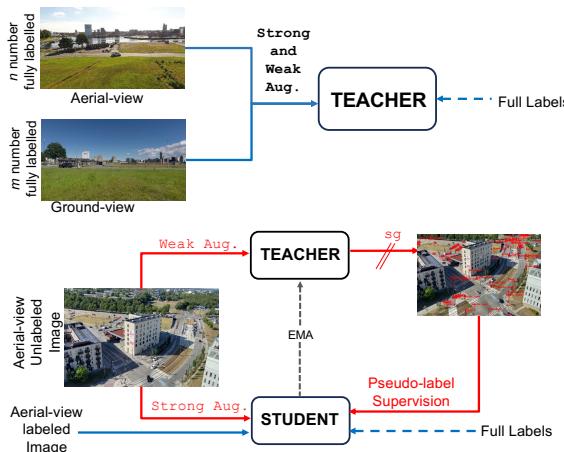


Figure 5. Semi-Supervised object detection framework based on curriculum learning approach. Here blue represents the initial supervised stage, red represents the later unsupervised stage, SG represents stop gradient.

shows the superiority of these architectures when large-scale data is available. (ii) These transformer based object detectors augmented with even 25% of MAVREC’s ground view images, surpasses the model pretrained on MSCOCO. This shows the importance of learning geography-aware representation for aerial visual perception, suggesting a new direction for enhancing object detection in aerial perspective.

**Model generalization.** In Table 3, we provide the object detection performance of MAVREC with models pretrained with different strategies. Our empirical evaluations indicate that state-of-the-art object detection models, including the *open-world foundational models* like Grounding-Dino [46], which fail to achieve the expected performance level on

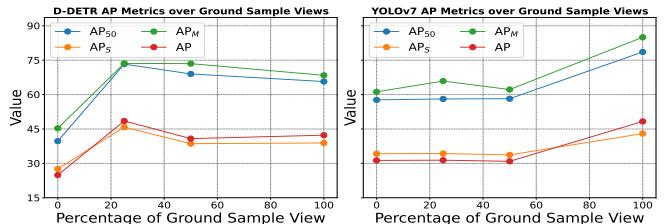


Figure 4. Supervised benchmark on aerial view of MAVREC (Validation Set).

MAVREC; see Figure 6. This observation validates an inherent bias of these models towards ground-view data. Moreover, Table 3, combined with the visual evidence in Figure 6, shows an object detector pre-trained on popular ground-view dataset (MS-COCO [44]) or other aerial datasets collected from different geographies (e.g., Visdrone [88] from China) has diminished efficacy on aerial images obtained from disparate geographical regions (for our case, Europe). Therefore, unlike classical object detection, training a sophisticated DNN model on a large dataset (e.g., ImageNet [64] or MS-COCO [44]) does not offer the best overall solution. We find augmenting object detectors with ground-view images from the corresponding geographical context is a superior strategy that boosts detection performance.

## 4.2. Curriculum learning based semi-supervised object detection

In this section, we introduce a curriculum based learning strategy for semi-supervised object detection. Curriculum learning [66] provides a systematic strategy to enhance model performance by incrementally introducing complexity into the training regime. We adopt curriculum learning in the semi-supervised object detection using a D-DETR. In our semi-supervised baseline, we train the object detector using both labeled and unlabeled image sets. The foundation of our semi-supervised baseline is a consistency regularization framework based on a teacher-student model. This framework encompasses two distinct training phases: (i) the ‘burn-in’ stage, where the teacher model is trained exclusively with labeled images, and (ii) the semi-supervised training stage, where the teacher-student model engages with unlabeled images. Given a trained teacher network from the burn-in stage, this framework leverages

464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495

Table 4. Semi-supervised Omni-DETR [74] benchmark on MAVREC. In the table,  $G$  and  $A$  denote number of ground and aerial-view images, respectively. During the burn-in, we only use the labelled subset.

Training Technique	Labelled		Unlabelled		Test perspective	Validation Set				Test Set			
	$G(m)$	$A(n)$	$G$	$A$		AP	AP <sub>50</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP	AP <sub>50</sub>	AP <sub>S</sub>	AP <sub>M</sub>
Semi-Supervised	8605	0	8605	0	G	56.9	83.3	54.8	74.9	45.8	75.5	45.4	58.7
Semi-Supervised	0	8605	0	8605	A	29.3	49.3	24.8	60.6	19.8	38.4	19.5	35.0
<b>Curriculum Learning (Ours)</b>	2151	8605	0	8605	A	37.8	<b>64.9</b>	35.4	68.3	23.2	45.4	21.8	43.6
<b>Curriculum Learning (Ours)</b>	2151	8605	2151	8605	A	<b>38.0</b>	64.8	<b>35.7</b>	67.6	<b>26.7</b>	<b>54.1</b>	<b>24.5</b>	42.4



Figure 6. **Left to right:** Sample frames from time-synchronized MAVREC showing aerial and ground views; D-DETR [89] trained on aerial VisDrone DET [88] inference results on MAVREC (GT bounding boxes are green, and detection are in red); Grouding-Dino [34] inference result on MAVERC; inference results of D-DETR trained on aerial MAVREC has fewer missed detections. (Zoom in for a better view)

weak-to-strong consistency regularization [22] to leverage unlabeled aerial images as shown in Figure 5. In this stage, the teacher network processes the weakly augmented unlabeled aerial-view images to generate bounding boxes. These bounding boxes serve as pseudo-labels for the strongly augmented counterpart of the image, which is the input of the student. The teacher network is updated through an exponential moving average (EMA) of the student’s updates, and the student network only gets back-propagated.

However, the performance of the above semi-supervised baseline partially relies on the effectiveness of the teacher network to generate pseudo-labels. Inspired from our insights in transferring geography-aware knowledge from ground view to aerial view, we employ a curriculum learning strategy in the burn-in stage of the semi-supervised framework. In contrast to training the teacher network with only labelled aerial images, we train the teacher, first, with  $m$  labelled ground-view images and then with  $n$  labelled aerial images. The outcome is a trained geography-aware teacher network that facilitates the second phase of training the semi-supervised framework by generating precise object proposals.

In Table 4, we showcase the results of semi-supervised benchmarking on the MAVREC dataset, utilizing D-DETR as the backbone object detection model within the teacher-student framework. Note that all the models are trained on 39 epochs with 20 epochs for the burn-in stage and 19 epochs for the semi-supervised training stage. Our semi-supervised baseline results demonstrate that by utilizing the same number of unlabeled aerial images as labeled images, we achieve a substantial boost in object detection performance—from 13.1% to 29.3% and 10.3% to 19.8% on validation and test set, respectively. We observe a consistent improvement in the ground view. This warrants the importance of using semi-

supervised approaches for vision tasks, particularly where the annotation process is labor-intensive. Subsequently, we demonstrate the performance of our curriculum-based semi-supervised object detector on the MAVREC dataset. This approach outperforms the baseline model by 8.5% and 7% on validation and test set, respectively. Moreover, we observed additional improvements in object detection accuracy when the second phase of training is augmented with unlabelled ground-view images. This shows the importance of learning geography-aware representation from the ground-view for enhancing aerial visual perception.

## 5. Conclusion

In this paper, we introduce a large-scale, high-definition ground and aerial-view video dataset, MAVREC. To the best of our knowledge, MAVREC is the first drone-based aerial object detection dataset that exploits the multi-view of the data coming from orthogonal views, aerial and ground to offer enhanced detection capacity for an aerial view. In our extensive benchmarking on MAVREC, we employed both supervised and semi-supervised learning methods, along with our proposed curriculum-based ground-view pre-training strategy. Our findings highlight several key insights, particularly the importance of geographic awareness in aerial visual models. We discovered that models trained on aerial images from one geographic location often struggle to generalize to different regions. However, integrating ground-view data from the same geographic area significantly enhances the model’s ability to learn more distinctive visual representations. We envision that this dataset and benchmarking will benefit: (i) researchers, who will use it as the basis for consistent implementation and evaluation; and (ii) practitioners, who need an appropriate, large-scale, industry-standard dataset for training DNN models for aerial images.

563 **References**

- [1] Color difference. [https://en.wikipedia.org/wiki/Color\\_difference](https://en.wikipedia.org/wiki/Color_difference). 16
- [2] CVAT annotation tool. <https://www.cvcat.ai>. 5, 15, 17
- [3] PNNL Parking Lot 1 and 2 and Pizza sequences. <https://www.crcv.ucf.edu/data/ParkingLOT/>. 14
- [4] The Most Popular Car Color: Can You Guess Which One?, . <https://www.motorbiscuit.com/most-popular-car-color-guess-color/>. 6
- [5] Most popular car-colors by country, . <https://haynes.com/en-us/tips-tutorials/most-popular-car-colors-country-or-don-t-buy-black-car-india>. 2
- [6] Innovation built through partnerships to improve life on the streetscape for all. <https://cs3-erc.org/>. 2, 22
- [7] DJI. <https://www.dji.com>. 3
- [8] Mapped: The World's Population Density by Latitude. <https://www.visualcapitalist.com/cp/mapped-the-worlds-population-density-by-latitude/>. 2, 15
- [9] Car colour popularity. [https://en.wikipedia.org/wiki/Car\\_colour\\_popularity](https://en.wikipedia.org/wiki/Car_colour_popularity), . 15
- [10] Solar zenith-angle, . [https://en.wikipedia.org/wiki/Solar\\_zenith\\_angle](https://en.wikipedia.org/wiki/Solar_zenith_angle). 14
- [11] Yolo-NAS. <https://github.com/Deci-AI/super-gradients/blob/master/YOLONAS.md>. 4, 6, 20, 21
- [12] Amirsaman Ashtari, Raehyuk Jung, Mingxiao Li, and Junyong Noh. A drone video clip dataset and its applications in automated cinematography. In *Computer Graphics Forum*, pages 189–203, 2022. 3
- [13] Naeem Ayoub and Peter Schneider-Kamp. Real-time on-board detection of components and faults in an autonomous uav system for power line inspection. In *Proceedings of the International Conference on Deep Learning Theory and Applications*, pages 68–75, 2020. 2, 22
- [14] Mohammadamin Barekatain, Miquel Martí, Hsueh-Fu Shih, Samuel Murray, Kotaro Nakayama, Yutaka Matsuo, and Helmut Prendinger. Okutama-action: An aerial view video dataset for concurrent human action detection. In *Proceedings of the Conference on computer vision and pattern recognition workshops*, pages 28–35, 2017. 14
- [15] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. End to end learning for self-driving cars. 2016. 2, 22
- [16] Elizabeth Bondi, Raghav Jain, Palash Aggrawal, Saket Anand, Robert Hannaford, Ashish Kapoor, Jim Piavis, Shital Shah, Lucas Joppa, Bistra Dilkina, et al. BIRD-SAI: A dataset for detection and tracking in aerial thermal infrared videos. In *Proceedings of the Winter Conference on Applications of Computer Vision*, pages 1747–1756, 2020. 2, 14
- [17] Ilker Bozcan and Erdal Kayacan. Au-air: A multi-modal unmanned aerial vehicle dataset for low altitude traffic surveillance. In *IEEE International Conference on Robotics and Automation*, pages 8504–8510, 2020. 2, 3
- [18] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Proceedings of the European Conference on Computer Vision*, pages 213–229, 2020. 4, 6, 17, 20, 21
- [19] Binghui Chen, Pengyu Li, Xiang Chen, Biao Wang, Lei Zhang, and Xian-Sheng Hua. Dense learning based semi-supervised object detection. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*, pages 4815–4824, 2022. 5
- [20] Jian Ding, Nan Xue, Gui-Song Xia, Xiang Bai, Wen Yang, Michael Ying Yang, Serge Belongie, Jiebo Luo, Mihai Datcu, Marcello Pelillo, et al. Object detection in aerial images: A large-scale benchmark and challenges. *IEEE transactions on pattern analysis and machine intelligence*, 44(11):7778–7796, 2021. 3
- [21] Dawei Du, Yuankai Qi, Hongyang Yu, Yifan Yang, Kaiwen Duan, Guorong Li, Weigang Zhang, Qingming Huang, and Qi Tian. The unmanned aerial vehicle benchmark: Object detection and tracking. In *Proceedings of the European Conference on Computer Vision*, pages 370–386, 2018. 2, 3
- [22] Yue Fan, Anna Kukleva, Dengxin Dai, and Bernt Schiele. Revisiting consistency regularization for semi-supervised learning. *International Journal of Computer Vision*, 131(3):626–643, 2023. 8
- [23] Dapeng Feng, Yuhua Qi, Shipeng Zhong, Zhiqiang Chen, Yudu Jiao, Qiming Chen, Tao Jiang, and Hongbo Chen. S3e: A large-scale multimodal dataset for collaborative slam. 2022. 3
- [24] Zheng Ge, Songtao Liu, Feng Wang, Zeming Li, and Jian Sun. Yolox: Exceeding yolo series in 2021. *arXiv preprint arXiv:2107.08430*, 2021. 4
- [25] Ross Girshick. Fast R-CNN. In *Proceedings of the International Conference on Computer Vision*, pages 1440–1448, 2015. 3
- [26] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Pro-*

- 668        *ceedings of the Conference on Computer Vision and*  
669        *Pattern Recognition*, pages 580–587, 2014. 3
- 670 [27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian  
671 Sun. Deep residual learning for image recognition. In  
672 *Proceedings of the Conference on Computer Vision and*  
673 *Pattern Recognition*, pages 770–778, 2016. 21
- 674 [28] Yu Hongyang, Guorong Li, Weigang Zhang, Qingming  
675 Huang, Dawei Du, Tian Qi, and Sebe Nicu. The un-  
676 manned aerial vehicle benchmark: Object detection,  
677 tracking and baseline. *International Journal of Com-*  
678 *puter Vision*, 128(5):1141–1159, 2020. 15
- 679 [29] Meng-Ru Hsieh, Yen-Liang Lin, and Winston H Hsu.  
680 Drone-based object counting by spatially regularized  
681 regional proposal network. In *Proceedings of the Inter-*  
682 *national Conference on Computer Vision*, pages 4145–  
683 4153, 2017. 14
- 684 [30] Zubayer Islam, Mohamed Abdel-Aty, Amrita  
685 Goswamy, Amr Abdelraouf, and Ou Zheng. Effect of  
686 signal timing on vehicles’ near misses at intersections.  
687 *Scientific reports*, 13:9065, 2023. 2, 22
- 688 [31] Efstratios Kakaletsis, Charalampos Symeonidis, Maria  
689 Tzelepi, Ioannis Mademlis, Anastasios Tefas, Nikos  
690 Nikolaidis, and Ioannis Pitas. Computer vision for au-  
691 tonomous UAV flight safety: an overview and a vision-  
692 based safe landing pipeline example. *ACM Computing*  
693 *Surveys*, 54(9):1–37, 2021. 2
- 694 [32] Isha Kalra, Maneet Singh, Shruti Nagpal, Richa Singh,  
695 Mayank Vatsa, and P. B. Sujit. DroneSURF: Bench-  
696 mark Dataset for Drone-based Face Recognition. In  
697 *proceedings of International Conference on Automatic*  
698 *Face and Gesture Recognition*, pages 1–7, 2019. 14
- 699 [33] Diederik P Kingma and Jimmy Ba. Adam: A method  
700 for stochastic optimization. In *Proceedings of the Inter-*  
701 *national Conference on Learning Representations*,  
702 2015. 20
- 703 [34] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi  
704 Mao, Chloe Rolland, Laura Gustafson, Tete Xiao,  
705 Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo,  
706 Piotr Dollár, and Ross Girshick. Segment anything.  
707 *arXiv:2304.02643*, 2023. 8
- 708 [35] Robert Krajewski, Julian Bock, Laurent Kloeker, and  
709 Lutz Eckstein. The hghd dataset: A drone dataset  
710 of naturalistic vehicle trajectories on german high-  
711 ways for validation of highly automated driving sys-  
712 tems. In *2018 21st International Conference on Intel-*  
713 *ligent Transportation Systems (ITSC)*, pages 2118–2125,  
714 2018. 14
- 715 [36] S.V. Aruna Kumar, Ehsan Yaghoubi, Abhijit Das, B.S.  
716 Harish, and Hugo Proen  a. The P-DESTRE: A fully  
717 annotated dataset for pedestrian detection, tracking,  
718 and short/long-term re-identification from aerial de-  
719 vices. *IEEE Transactions on Information Forensics*  
720 *and Security*, 16:1696–1708, 2020. 15
- [37] Matti Kummu and Olli Varis. The world by latitudes:  
721 A global analysis of human population, development  
722 level and environment across the north–south axis over  
723 the past half century. *Applied geography*, 31(2):495–  
724 507, 2011. 2, 15
- [38] Matti Kummu, Hans De Moel, Gianluigi Salvucci,  
725 Daniel Viviroli, Philip J Ward, and Olli Varis. Over  
726 the hills and further away from coast: global geospatial  
727 patterns of human and environment over the 20th–  
728 21st centuries. *Environmental Research Letters*, 11(3):  
729 034010, 2016. 2
- [39] Aoxue Li, Peng Yuan, and Zhenguo Li. Semi-  
730 supervised object detection via multi-instance align-  
731 ment with global class prototypes. In *Proceedings of*  
732 *the Conference on Computer Vision and Pattern Recog-*  
733 *nition*, pages 9809–9818, 2022. 5
- [40] Chuyi Li, Lulu Li, Hongliang Jiang, Kaiheng Weng,  
734 Yifei Geng, Liang Li, Zaidan Ke, Qingyuan Li, Meng  
735 Cheng, Weiqiang Nie, et al. Yolov6: A single-stage  
736 object detection framework for industrial applications.  
*arXiv preprint arXiv:2209.02976*, 2022. 4
- [41] Jing Li, Dong Hye Ye, Timothy Chung, Mathias  
737 Kolsch, Juan Wachs, and Charles Bouman. Multi-  
738 target detection and tracking from a single camera in  
739 unmanned aerial vehicles (UAVs). In *Proceedings of*  
740 *the International Conference on Intelligent Robots and*  
741 *Systems*, pages 4992–4997, 2016. 14
- [42] Jingtong Li, Jesse Murray, Dorina Ismaili, Konrad  
742 Schindler, and Cenek Albl. Reconstruction of 3D flight  
743 trajectories from ad-hoc camera networks. In *Pro-*  
744 *ceedings of the International Conference on Intelligent*  
745 *Robots and Systems*, pages 1621–1628, 2020. 14
- [43] Siyi Li and Dit-Yan Yeung. Visual object tracking  
746 for unmanned aerial vehicles: A benchmark and new  
747 motion models. In *Proceedings of the AAAI Conference*  
748 *on Artificial Intelligence*, 2017. 15
- [44] Tsung-Yi Lin, Michael Maire, Serge Belongie, James  
749 Hays, Pietro Perona, Deva Ramanan, Piotr Doll  r, and  
750 C Lawrence Zitnick. Microsoft COCO: Common ob-  
751 jects in context. In *Proceedings of the European Con-*  
752 *ference on Computer Vision*, pages 740–755, 2014. 1,  
753 5, 6, 7, 20
- [45] Tsung-Yi Lin, Yin Cui, Serge Belongie, and James  
754 Hays. Learning deep representations for ground-to-  
755 aerial geolocation. In *Proceedings of the Con-*  
756 *ference on Computer Vision and Pattern Recognition*,  
757 pages 5007–5015, 2015. 14
- [46] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li,  
758 Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang,  
759 Hang Su, Jun Zhu, et al. Grounding dino: Marrying  
760 dino with grounded pre-training for open-set object  
761 detection. *arXiv preprint arXiv:2303.05499*, 2023. 7

- 773 [47] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Chris- 825  
774 tian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander 826  
775 Berg. SSD: Single shot multibox detector. In 827  
776 *Proceedings of the European Conference on Computer 828  
777 Vision*, pages 21–37, 2016. 4  
778 [48] Yen-Cheng Liu, Chih-Yao Ma, and Zsolt Kira. Unbi- 829  
779 ased teacher v2: Semi-supervised object detection for 830  
780 anchor-free and anchor-based detectors. In *Proceed- 831  
781 ings of the Conference on Computer Vision and Pattern 832  
782 Recognition*, pages 9819–9828, 2022. 5  
783 [49] Ye Lyu, George Vosselman, Gui-Song Xia, Alper Yil- 833  
784 maz, and Michael Ying Yang. UAVid: A semantic 834  
785 segmentation dataset for uav imagery. *ISPRS journal 835  
786 of photogrammetry and remote sensing*, 165:108–119, 836  
787 2020. 2, 3  
788 [50] András L Majdik, Damiano Verda, Yves Albers- 837  
789 Schoenberg, and Davide Scaramuzza. Air-ground 838  
790 matching: Appearance-based GPS-denied urban 839  
791 localization of micro aerial vehicles. *Journal of Field 840  
792 Robotics*, 32(7):1015–1039, 2015. 14  
793 [51] Murari Mandal, Lav Kush Kumar, and Santosh Ku- 841  
794 mar Vipparthi. Mor-UAV: A benchmark dataset and 842  
795 baselines for moving object recognition in UAV videos. 843  
796 In *Proceedings of ACM International Conference on 844  
797 Multimedia*, pages 2626–2635, 2020. 2, 3  
798 [52] Aboli Marathe, Pushkar Jain, Rahee Walambe, and 845  
799 Ketan Kotecha. RestoreX-AI: A contrastive approach 846  
800 towards guiding image restoration via explainable AI 847  
801 systems. In *Proceedings of the Conference on Com- 848  
802 puter Vision and Pattern Recognition Workshops*, pages 849  
803 3030–3039, 2022. 2, 22  
804 [53] Peng Mi, Jianghang Lin, Yiyi Zhou, Yunhang Shen, 850  
805 Gen Luo, Xiaoshuai Sun, Liujuan Cao, Rongrong Fu, 851  
806 Qiang Xu, and Rongrong Ji. Active teacher for semi- 852  
807 supervised object detection. In *Proceedings of the Con- 853  
808 ference on Computer Vision and Pattern Recognition*, 854  
809 pages 14482–14491, 2022. 5  
810 [54] Matthias Mueller, Neil Smith, and Bernard Ghanem. 855  
811 A benchmark and simulator for UAV tracking. In *Pro- 856  
812 ceedings of the European Conference on Computer 857  
813 Vision*, pages 445–461, 2016. 2, 3, 17, 18, 19  
814 [55] Arjun Nagendran, Don Harper, and Mubarak Shah. 858  
815 New system performs persistent wide-area aerial 859  
816 surveillance. *SPIE Newsroom*, 5:20–28, 2010. 14  
817 [56] Priya Narayanan, Xin Hu, Zhenyu Wu, Matthew D. 860  
818 Thielke, John G. Rogers, Andre V Harrison, John A. 861  
819 D’Agostino, James D Brown, Long P. Quang, James R. 862  
820 Uplinger, Heesung Kwon, and Zhangyang Wang. A 863  
821 multi-purpose realistic haze benchmark with quantifi- 864  
822 able haze levels and ground truth. *IEEE Transactions 865  
823 on Image Processing*, 32:3481–3492, 2023. 3  
824 [57] Yurii Nesterov. *Introductory lectures on convex opti- 866  
825 mization: A basic course*. Springer Science and Busi- 867  
826 ness Media, 2003. 20  
827 [58] Sangmin Oh, Anthony Hoogs, Amitha Perera, Naresh 828  
828 Cuntoor, Chia-Chih Chen, Jong Taek Lee, Saurajit 829  
829 Mukherjee, J. K. Aggarwal, Hyungtae Lee, Larry 830  
830 Davis, Eran Swears, Xioyang Wang, Qiang Ji, Kishore 831  
831 Reddy, Mubarak Shah, Carl Vondrick, Hamed Pirsi- 832  
832 avash, Deva Ramanan, Jenny Yuen, Antonio Torralba, 833  
833 Bi Song, Anesco Fong, Amit Roy-Chowdhury, and 834  
834 Mita Desai. A large-scale benchmark dataset for event 835  
835 recognition in surveillance video. In *Proceedings of 836  
836 the Conference on Computer Vision and Pattern Rec- 837  
837ognition*, pages 3153–3160, 2011. 15  
838 [59] Yaoru Pan, Mogens Flindt, Peter Schneider-Kamp, and 838  
839 Marianne Holmer. Beach wrack mapping using un- 840  
840 manned aerial vehicles for coastal environmental man- 841  
841 agement. *Ocean and Coastal Management*, 213, 2021. 842  
842 2, 22  
843 [60] Anne-Flore Perrin, Vassilios Krassanakis, Lu Zhang, 843  
844 Vincent Ricordel, Matthieu Perreira Da Silva, and 844  
845 Olivier Le Meur. EyetrackUAV2: A large-scale binoc- 845  
846 ular eye-tracking dataset for UAV videos. *Drones*, 4 846  
847 (1):2, 2020. 14  
848 [61] N Dinesh Reddy, Minh Vo, and Srinivasa G 848  
849 Narasimhan. Carfusion: Combining point tracking and 849  
850 part detection for dynamic 3D reconstruction of ve- 850  
851 hicles. In *Proceedings of the Conference on Computer 851  
852 Vision and Pattern Recognition*, pages 1906–1915, 2018. 852  
853 3, 14  
854 [62] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian 854  
855 Sun. Faster R-CNN: Towards real-time object detec- 855  
856 tion with region proposal networks. *Proceedings of 856  
857 Advances in neural information processing systems*, 28, 857  
858 2015. 3  
859 [63] Alexandre Robicquet, Amir Sadeghian, Alexandre 859  
860 Alahi, and Silvio Savarese. Learning social etiquette: 860  
861 Human trajectory understanding in crowded scenes. In 861  
862 *Proceedings of the European Conference on Computer 862  
863 Vision*, pages 549–565, 2016. 2, 3, 17, 18, 19  
864 [64] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, 864  
865 Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej 865  
866 Karpathy, Aditya Khosla, Michael Bernstein, Alexan- 866  
867 der C. Berg, and Li Fei-Fei. Imagenet large scale visual 867  
868 recognition challenge, 2014. 7  
869 [65] Qi Shan, Changchang Wu, Brian Curless, Yasutaka Fu- 869  
870 rukawa, Carlos Hernandez, and Steven M Seitz. Accu- 870  
871 rate geo-registration by ground-to-aerial image match- 871  
872 ing. In *Proceedings of International Conference on 3D 872  
873 Vision*, pages 525–532, 2014. 14  
874 [66] Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and 874  
875 Nicu Sebe. Curriculum learning: A survey. *Interna- 875  
876 tional Journal of Computer Vision*, 130(6):1526–1565, 876  
877 2022. 7

- 878 [67] Yiming Sun, Bing Cao, Pengfei Zhu, and Qinghua  
879 Hu. Drone-based RGB-infrared cross-modality vehicle  
880 detection via uncertainty-aware learning. *IEEE Transactions  
881 on Circuits and Systems for Video Technology*,  
882 32(10):6700–6713, 2022. 2, 14
- 883 [68] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He.  
884 FCOS: A simple and strong anchor-free object detector.  
885 *IEEE Transactions on Pattern Analysis and Machine  
886 Intelligence*, 44(4):1922–1933, 2020. 4
- 887 [69] Andrea Vallone, Frederik Warburg, Hans Hansen,  
888 Søren Hauberg, and Javier Civera. Danish airs and  
889 grounds: A dataset for aerial-to-street-level place recog-  
890 nition and localization. *IEEE Robotics and Automation  
891 Letters*, 7(4):9207–9214, 2022. 14
- 892 [70] Leon Amadeus Varga, Benjamin Kiefer, Martin Mess-  
893 mer, and Andreas Zell. Seadroneesee: A maritime  
894 benchmark for detecting humans in open water. In  
895 *Proceedings of the Winter Conference on Applications  
896 of Computer Vision*, pages 2260–2270, 2022. 14
- 897 [71] Chuanyun Wang, Yang Su, Jingjing Wang, Tian Wang,  
898 and Qian Gao. UAVSwarm dataset: An unmanned  
899 aerial vehicle swarm dataset for multiple object track-  
900 ing. *Remote Sensing*, 14(11), 2022. 14
- 901 [72] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-  
902 Yuan Mark Liao. Scaled-yolov4: Scaling cross stage  
903 partial network. In *Proceedings of the Conference  
904 on Computer Vision and Pattern Recognition*, pages  
905 13029–13038, 2021. 4
- 906 [73] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-  
907 Yuan Mark Liao. YoloV7: Trainable bag-of-freebies  
908 sets new state-of-the-art for real-time object detectors.  
909 In *Proceedings of the Conference on Computer Vision  
910 and Pattern Recognition*, pages 7464–7475, 2023. 4, 6,  
911 20, 21
- 912 [74] Pei Wang, Zhaowei Cai, Hao Yang, Gurumurthy  
913 Swaminathan, Nuno Vasconcelos, Bernt Schiele, and  
914 Stefano Soatto. Omni-DETR: Omni-supervised object  
915 detection with transformers. In *Proceedings of the Con-  
916 ference on Computer Vision and Pattern Recognition*,  
917 pages 9367–9376, 2022. 5, 8, 20, 21
- 918 [75] Xin Wang, Ning He, Chen Hong, Qi Wang, and Ming  
919 Chen. Improved YOLOX-X based UAV aerial photog-  
920 raphy object detection algorithm. *Image and Vision  
921 Computing*, 135:104697, 2023. 4
- 922 [76] Xinjiang Wang, Xingyi Yang, Shilong Zhang, Yijiang  
923 Li, Litong Feng, Shijie Fang, Chengqi Lyu, Kai Chen,  
924 and Wayne Zhang. Consistent-Teacher: Towards re-  
925 ducing inconsistent pseudo-targets in semi-supervised  
926 object detection. In *Proceedings of the Conference  
927 on Computer Vision and Pattern Recognition*, pages  
928 3240–3249, 2023. 5
- 929 [77] Longyin Wen, Dawei Du, Pengfei Zhu, Qinghua Hu,  
930 Qilong Wang, Liefeng Bo, and Siwei Lyu. Detec-  
tion, tracking, and counting meets drones in crowds:  
931 A benchmark. In *Proceedings of the Conference on  
932 Computer Vision and Pattern Recognition*, pages 7808–  
933 7817, 2021. 15
- 934 [78] Xin Wu, Wei Li, Danfeng Hong, Ran Tao, and Qian  
935 Du. Deep learning for unmanned aerial vehicle-based  
936 object detection and tracking: A survey. *IEEE Geo-  
937 science and Remote Sensing Magazine*, 10(1):91–124,  
938 2022. 2, 3, 14, 15, 22
- 939 [79] Gui-Song Xia, Xiang Bai, Jian Ding, Zhen Zhu, Serge  
940 Belongie, Jiebo Luo, Mihai Datcu, Marcello Pelillo,  
941 and Liangpei Zhang. DOTA: A large-scale dataset  
942 for object detection in aerial images. In *Proceedings  
943 of the Conference on Computer Vision and Pattern  
944 Recognition*, pages 3974–3983, 2018. 3
- 945 [80] Chang Xu, Jinwang Wang, Wen Yang, Huai Yu, Lei  
946 Yu, and Gui-Song Xia. Detecting tiny objects in aerial  
947 images: A normalized wasserstein distance and a new  
948 benchmark. *ISPRS Journal of Photogrammetry and  
949 Remote Sensing*, 190:79–93, 2022. 4
- 950 [81] Chang Xu, Jinwang Wang, Wen Yang, Huai Yu, Lei  
951 Yu, and Gui-Song Xia. RFLA: Gaussian receptive field  
952 based label assignment for tiny object detection. In  
953 *Proceedings of the European Conference on Computer  
954 Vision*, pages 526–543, 2022. 4
- 955 [82] Xiaowei Xu, Xinyi Zhang, Bei Yu, Xiaobo Sharon Hu,  
956 Christopher Rowen, Jingtong Hu, and Yiyu Shi. DAC-  
957 SDC low power object detection challenge for UAV  
958 applications. *IEEE Transactions on Pattern Analysis  
959 and Machine Intelligence*, 43(2):392–403, 2021. 3
- 960 [83] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin,  
961 and David Lopez-Paz. Mixup: Beyond Empirical Risk  
962 Minimization. In *Proceedings of the International  
963 Conference on Learning Representations*, 2018. 21
- 964 [84] Haijun Zhang, Mingshan Sun, Qun Li, Linlin Liu,  
965 Ming Liu, and Yuzhu Ji. An empirical study of multi-  
966 scale object detection in high resolution UAV images.  
967 *Neurocomputing*, 421:173–182, 2021. 3
- 968 [85] Wang Zhang, Chunsheng Liu, Faliang Chang, and Ye  
969 Song. Multi-scale and occlusion aware network for ve-  
970 hicle detection and segmentation on uav aerial images.  
971 *Remote Sensing*, 12(11):1760, 2020. 14
- 972 [86] Ou Zheng, Mohamed Abdel-Aty, Lishengsa Yue, Amr  
973 Abdelraouf, Zijin Wang, and Nada Mahmoud. CitySim:  
974 A drone-based vehicle trajectory dataset for safety  
975 oriented research and digital twins. *arXiv preprint  
976 arXiv:2208.11036*, 2022. 2, 22
- 977 [87] Pengfei Zhu, Jiayu Zheng, Dawei Du, Longyin Wen,  
978 Yiming Sun, and Qinghua Hu. Multi-drone-based sin-  
979 gle object tracking with agent sharing network. *IEEE  
980 Transactions on Circuits and Systems for Video Tech-  
981 nology*, 31(10):4058–4070, 2020. 2, 3

- 983 [88] Pengfei Zhu, Longyin Wen, Dawei Du, Xiao Bian,  
984 Heng Fan, Qinghua Hu, and Haibin Ling. Detection  
985 and tracking meet drones challenge. *IEEE Transactions*  
986 on Pattern Analysis and Machine Intelligence
- 987 , 44(11):7380–7399, 2022. 1, 2, 3, 5, 7, 8, 17, 18, 19
- 988 [89] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang  
989 Wang, and Jifeng Dai. Deformable DETR: Deformable  
990 Transformers for End-to-End Object Detection. In *Pro-*  
991 *ceedings of the International Conference on Learning*  
992 *Representations*, 2020. 4, 6, 7, 8, 17, 20, 21
- 993 [90] Xingkui Zhu, Shuchang Lyu, Xu Wang, and Qi Zhao.  
994 TPH-YOLOv5: Improved YOLOv5 based on trans-  
995 former prediction head for object detection on drone-  
996 captured scenarios. In *Proceedings of the International*  
997 *Conference on Computer Vision*, pages 2778–2788,  
998 2021. 4
- 999 [91] Yilin Zhu, Yang Kong, Yingrui Jie, Shiyou Xu, and  
1000 Hui Cheng. Graco: A multimodal dataset for ground  
1001 and aerial cooperative localization and mapping. *IEEE*  
1002 *Robotics and Automation Letters*, 2023. 3