

# Crowd Detection Using Very-Fine-Resolution Satellite Imagery

Tong Xiao <sup>a, b</sup>, Qunming Wang <sup>b, \*</sup>, Ping Lu <sup>b</sup>, Tenghai Huang <sup>b</sup>, Xiaohua Tong <sup>b</sup>, Peter M. Atkinson <sup>c, d</sup>

<sup>a</sup> Key Laboratory of Ethnic Language Intelligent Analysis and Security Governance of MOE, Minzu University of China, Beijing 100081, China

<sup>b</sup> College of Surveying and Geo-Informatics, Tongji University, 1239 Siping Road, Shanghai 200092, China

<sup>c</sup> Faculty of Science and Technology, Lancaster University, Lancaster LA1 4YR, UK

<sup>d</sup> Geography and Environment, University of Southampton, Highfield, Southampton SO17 1BJ, UK

\*Corresponding author. Email: wqm11111@126.com

**Abstract:** Accurate crowd detection (CD) is critical for public safety and historical pattern analysis, yet existing methods relying on ground and aerial imagery suffer from limited spatio-temporal coverage. The development of very-fine-resolution (VFR) satellite sensor imagery (e.g., ~0.3 m spatial resolution) provides unprecedented opportunities for large-scale crowd activity analysis, but it has never been considered for this task. To address this gap, we proposed CrowdSat-Net, a novel point-based convolutional neural network, which features two innovative components: Dual-Context Progressive Attention Network (DCPAN) to improve feature representation of individuals by aggregating scene context and local individual characteristics, and High-Frequency Guided Deformable Upsampler (HFGDU) that recovers high-frequency information during upsampling through frequency-domain guided deformable convolutions. To validate the effectiveness of CrowdSat-Net, we developed CrowdSat, the first VFR satellite imagery dataset designed specifically for CD tasks, comprising over 120k manually labeled individuals from multi-source satellite platforms (Beijing-3N, Jilin-1 Gaofen-04A and Google Earth) across China. In the experiments, CrowdSat-Net was compared with eight state-of-the-art point-based CD methods (originally designed for ground or aerial imagery and satellite-based animal detection) using CrowdSat and achieved the largest F1-score of 66.12% and Precision of 73.23%, surpassing the second-best method by 1.71% and 2.42%, respectively. Moreover,

extensive ablation experiments validated the importance of the DCPAN and HFGDU modules. Furthermore, cross-regional evaluation further demonstrated the spatial generalizability of CrowdSat-Net. This research advances CD capability by providing both a newly developed network architecture for CD and a pioneering benchmark dataset to facilitate future CD development. The source code is available at <https://github.com/Tong-777777/CrowdSat-Net>.

**Keywords:** Crowd detection; Very-fine-resolution (VFR) satellite imagery; Deep learning; Feature enhancement; Feature fusion.

## 1. Introduction

In recent years, population growth, continued urbanization and rapid economic development in many regions have led to an increase in the frequency of crowd activities in public areas. Such crowds can pose a series of public safety risks, including traffic congestion (Gazzawe and Albahar, 2024), crowd crushes (Al-Nami, 2023), security incidents (Feliciani et al., 2022) and public health risks (Pokhrel and Chhetri, 2021; Joiner et al., 2024). Alleviating these risks, for example through crowd control, has become an essential aspect of ensuring public safety in densely populated gatherings. Crowd detection (CD), which involves estimating the location and count of the individuals in a crowd in a specific place (Sam et al., 2020; Wan et al., 2021; Han et al., 2023), plays a crucial role in mitigating these risks. By managing and controlling crowd movements (Weng et al., 2023; Rezaee et al., 2024) effectively, CD can help reduce the likelihood of accidents and create a safer environment. In addition to crowd management, long-term CD can analyze future crowd distribution patterns by mining historical data trends, providing valuable insights for urban planning and infrastructure optimization (Luo et al., 2025). The definition of a crowd varies across research domains. In this research, we

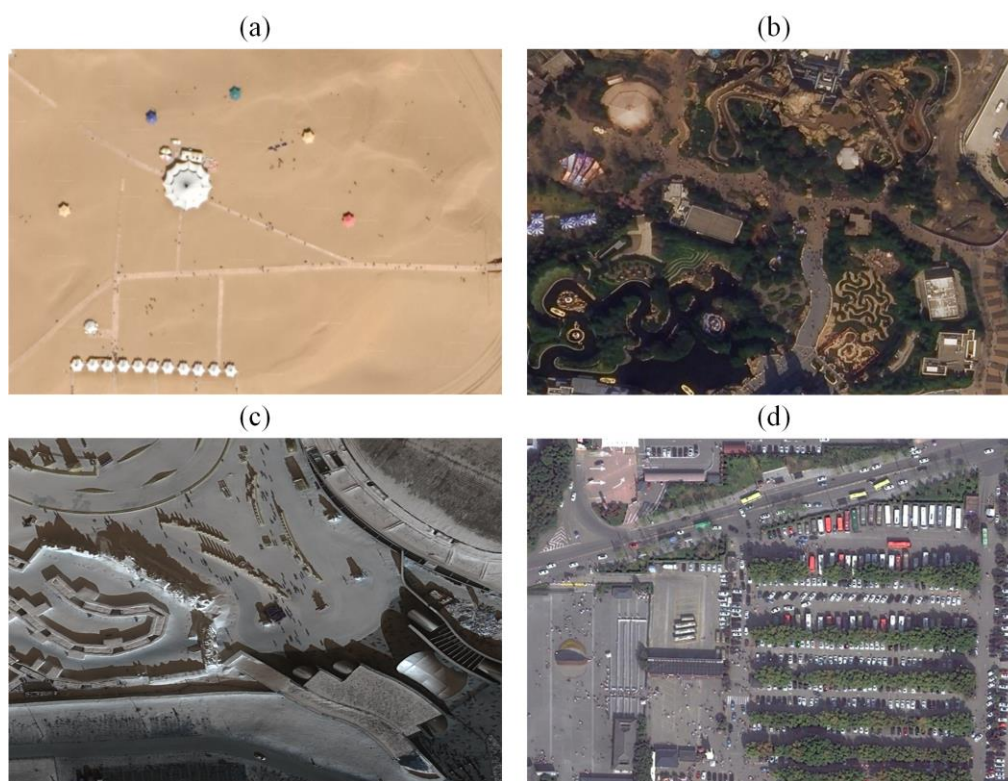
51 adopt a broad definition: A gathering of individuals, whether isolated, dispersed or clustered, or any  
52 combination of these, is regarded as a crowd for detection.

53 To achieve CD effectively, early research relied primarily on analyzing imagery captured by ground  
54 surveillance or fixed cameras (Lempitsky and Zisserman, 2010; Zhang et al., 2021; Mo et al., 2024). These  
55 cameras provide diverse perspectives, varying scales and different illumination conditions with which to  
56 detect crowds. However, their coverage was inherently limited to specific areas within the camera's field of  
57 view. With the development of airborne and unmanned aerial vehicle platforms, aerial datasets, such as the  
58 DLR Aerial Crowd Dataset (Bahmanyar et al., 2019), Meynberg's dataset (Meynberg et al., 2016) and Mliki's  
59 dataset (Mliki et al., 2019), have been introduced for CD, enhancing the development of models to monitor  
60 and analyze crowd dynamics over large areas. Although aerial imagery provides many advantages for  
61 monitoring, especially its potentially fine spatial resolution, its effectiveness for long-term analysis is  
62 constrained by infrequent data updates. For example, in Hebei Province, China, the standard update cycle for  
63 aerial photography typically exceeds three years (Hebei Provincial People's Government, 2011), which limits  
64 its ability to capture continuous short-term variations in crowd distributions.

65 With advances in remote sensing satellite technology, fine spatial-temporal resolution imagery has  
66 increasingly been applied to detect small-scale (i.e., small area) objects, such as ships, airplanes and tiny  
67 buildings (Gao et al., 2020a,b; Yi et al., 2023b). Compared to aerial imagery, satellite sensor imagery  
68 (hereafter, satellite imagery) offers a broader spatial coverage and shorter revisit intervals (e.g., the  
69 Worldview-4 satellite has a revisit period of approximately 4.5 days (Rumora et al., 2020)), increasing the  
70 potential for consistent, large-scale monitoring and historical analysis. Nowadays, very-fine-resolution (VFR)  
71 satellite imagery from sources like the Beijing-3N (BJ3N) and Jilin-1 Gaofen-04A (JL4A) satellites, as well as  
72 the Google Earth platform, has enabled the discrimination of individuals with the naked eye (as shown in Fig.  
73 1). This advance potentially allows more detailed and extensive monitoring, providing a valuable tool for  
74 applications requiring accurate object detection and tracking over wide areas. In light of this, as a primary  
75 contribution of this research, we presented a novel CD dataset, namely, CrowdSat, derived from the BJ3N and

76 JL4A satellites, as well as the Google Earth platform, covering 32 provincial-level divisions in China (except  
77 Guizhou Province and Macao). CrowdSat contains over 120k manually labeled individuals and consists of  
78 diverse regions with strong heterogeneity (e.g., built-up areas, snowy regions, beaches, desert regions, etc.).  
79 Compared to ground imagery, CrowdSat provides a much broader coverage, making it a more appropriate  
80 source for large-scale monitoring. Moreover, with its high revisit frequency, CrowdSat enables more  
81 continuous analysis of human activity patterns compared to aerial imagery, which is typically updated less  
82 frequently, especially at a large scale. To the best of our knowledge, CrowdSat is the first-ever CD dataset  
83 based on VFR satellite imagery, and it is intended to facilitate research on large-scale crowd analysis while  
84 uncovering historical patterns of human movement.

85



86

87 Fig. 1. Examples of VFR satellite imagery for CD. (a) Yuesha Island, Inner Mongolia, China, acquired by the Google Earth platform  
88 on Feb. 22, 2023 (© Google Earth 2025) with a spatial resolution of 0.30 m. (b) Shanghai Disney, China, acquired by the JL4A  
89 satellite on Feb. 16, 2023 with a spatial resolution of 0.31 m. (c) Harbin Ice and Snow World on Feb. 2, 2024 and (d) the Xi'an  
90 Emperor Qinshihuang's Mausoleum Site Museum in China on Oct. 29, 2023, both acquired by the BJ3N satellite with a spatial  
91 resolution of 0.30 m.

92

93 Compared to ground and aerial imagery, VFR satellite imagery is more stable and provided with greater  
94 consistency through time and can, thus, capture several sudden or unexpected crowd gathering cases. VFR  
95 satellite imagery, thus, has a greater range of potential applications. For example, VFR imagery can support  
96 post-disaster assessment and emergency management, long-term monitoring of urban mobility and public  
97 space usage, large-event planning and evaluation, and humanitarian and public-health interventions in refugee  
98 or quarantine areas. It can also assist with the sustainable management of tourism and cultural heritage sites by  
99 revealing how crowd densities evolve through time at popular destinations. These diverse applications  
100 highlight the need for advanced CD techniques based on VFR satellite imagery, particularly for tasks requiring  
101 precise individual-level information.

102 In the realm of CD, previous methods focused primarily on estimating crowd counts ([Shu et al., 2022](#); [Khan](#)  
103 [et al., 2023](#); [Yi et al., 2023a](#)). Counting, however, provides only coarse information and is insufficient for  
104 high-level tasks, such as individual tracking, activity recognition, anomaly detection and crowd flow or  
105 behavior prediction ([Laradji et al., 2018](#); [Song et al., 2021](#)). Recent research, therefore, pays more attention to  
106 fine-grained localization of individuals and formulates CD as segmentation- ([Meynberg et al., 2016](#)), density  
107 map- ([von Borstel et al., 2016](#)), detection- ([Ren et al., 2015](#)) or point-based prediction ([Zhou et al., 2019, 2020](#)).  
108 In this research, we adopt a point-based formulation, which represents each person by a single point while  
109 retaining accurate localization and requiring relatively simple manual annotation.

110 Notably, the aforementioned approaches are employed primarily in images captured by ground surveillance  
111 or fixed cameras, or aerial vehicles. As illustrated in [Fig. 2](#), individuals in such imagery exhibit typically clear  
112 and distinguishable features. When these clear signals become tiny and ambiguous in VFR satellite imagery, it  
113 is unclear whether the aforementioned methods can still work effectively. In the remote sensing field, related  
114 progress has already been made on dense small-object detection or counting using satellite imagery. [Yuan et al.](#)  
115 [\(2025\)](#) conducted comparison between YOLOv11, Single Shot MultiBox Detector (SSD), Faster  
116 Region-based Convolutional Neural Network (R-CNN) and Cascade R-CNN, for car detection using xView

117 images and bee-box detection using SkySat images, which shows that multi-scale feature extraction and  
118 anchor-free design are important for small objects. [Guo et al. \(2022\)](#) proposed a density map-based vehicle  
119 counting method for remote sensing imagery with limited spatial resolution, demonstrating that regressing  
120 vehicle density maps can yield accurate counts even when individual vehicles occupy only a few pixels. [Bashir  
121 and Wang \(2021\)](#) combined residual feature aggregation-based super-resolution with an object detector to  
122 enhance small-object detection in satellite and aerial images. [Delplanque et al. \(2023\)](#) proposed a single-stage  
123 CNN, which localizes each animal as a point. [Wu et al. \(2023\)](#) developed an ensemble of U-Net-based models  
124 to detect wildebeests from VFR satellite imagery over heterogeneous savanna landscapes, converting  
125 annotations to segmentation masks and using context-aware post-processing to obtain individual counts.  
126 These studies show that fine spatial resolution feature maps, multi-scale feature fusion and contextual priors  
127 can recover partially tiny and ambiguous signals of small objects in satellite imagery. Dense crowd detection  
128 using VFR satellite imagery, however, presents additional challenges. First, the signal of each individual in  
129 VFR satellite imagery is weaker than that of the small objects considered in previous remote sensing studies.  
130 As shown in Fig. 1, the signal size of each individual is approximately  $3 \times 3$  pixels. During the convolution  
131 process, particularly in the pooling stages, this small-sized signal can lead to attenuation or even loss of the  
132 signal ([Liu et al., 2021](#); [Tong and Wu, 2022](#); [Wei et al., 2024](#)). Second, multi-scale feature fusion, which  
133 integrates the upsampled coarse and fine spatial resolution features, is potentially a useful method to alleviate  
134 this issue. However, traditional upsampling methods, such as nearest neighbor and bilinear interpolation, often  
135 fail to recover fine spatial details, leading to over-smooth boundaries ([Li et al., 2024](#)) and misalignment  
136 between the high-frequency details in the upsampled features and fine spatial resolution features ([Dai et al.,  
137 2017](#)), making it challenging to detect small-sized objects effectively.

138



139

140 Fig. 2. Examples of ground and aerial imagery. (a) Image from the ShanghaiTech Dataset<sup>1</sup>, and (b) image from the DLR Aerial  
 141 Crowd Dataset<sup>2</sup>.

142

143 To overcome the aforementioned limitations, this paper proposed a novel point-based convolutional neural  
 144 network (CNN) method, CrowdSat-Net, which was specifically designed for large-scale and long-term CD.  
 145 CrowdSat-Net introduces two key contributions: 1) a Dual-Context Progressive Attention Network (DCPAN)  
 146 to improve the individual instance feature presentation and 2) a High-Frequency Guided Deformable  
 147 Upsampler (HFGDU) to replace traditional upsampling methods, which aims to recover the fine spatial  
 148 information of individuals during the upsampling process.

149 In summary, the contributions of this research are three-fold:

- 150 1) To the best of our knowledge, this is the first research to utilize VFR satellite imagery for CD, which  
 151 aims to facilitate studies on the characterization of human spatial distributions and temporal activities at

<sup>1</sup> <https://github.com/desenzhou/ShanghaiTechDataset>

<sup>2</sup> <https://www.dlr.de/en/eoc/about-us/remote-sensing-technology-institute/photogrammetry-and-image-analysis/public-datasets/dlr-acd>

152 a large-scale (both spatially and temporally).

- 153 2) To achieve this task, a novel CD dataset, CrowdSat, was collected by multi-source satellite platforms,  
154 which comprises over 120k labeled individuals and consists of diverse regions with strong  
155 heterogeneity, facilitating the development of CD methods. Additionally, during labeling, a point-like  
156 background removal strategy was introduced that uses multi-temporal VFR satellite imagery as  
157 auxiliary data to reduce mislabeling rates.
- 158 3) A novel point-based CD method using the CNN, termed CrowdSat-Net, was proposed to detect  
159 individuals in satellite imagery efficiently. Additionally, two innovational modules, DCPAN and  
160 HFGDU, were introduced to enhance individual feature presentation and recover the lost  
161 high-frequency information of individual features during upsampling.

162 The remainder of this paper is structured as follows: Section 2 provides a comprehensive overview of  
163 CrowdSat, detailing its data collection and preprocessing, labeling process and data analysis. Section 3  
164 presents the architecture of CrowdSat-Net, explaining its key components and design. Section 4 presents  
165 extensive experimental evaluations, while Section 5 discusses the broader applicability, inherent limitations of  
166 CrowdSat and CrowdSat-Net, and future directions. Finally, Section 6 summarizes the key findings of this  
167 research.

## 170 2. CrowdSat Dataset

171

172 In this paper, CrowdSat was presented for large-scale CD and the analysis of historical human movement  
173 patterns. To better understand the details of CrowdSat, this section provides a comprehensive overview,  
174 covering three key aspects: data collection and preprocessing, data labeling and data analysis.

### 176 2.1. Data Collection and Preprocessing

177

178 To ensure representative data for large-scale CD, three complementary remote sensing data sources,  
179 including the Google Earth platform and the BJ3N and JL04A satellites, were considered. The Google Earth  
180 imagery used in CrowdSat has an approximate spatial resolution of 0.30 m and originates primarily from  
181 Maxar (DigitalGlobe) WorldView-series satellites (e.g., WorldView-3/4). These data are pan-sharpened  
182 products produced by fusion of 0.30 m panchromatic and 1.20 m multispectral bands. It is noted that all  
183 Google Earth imagery used in this research was accessed solely for academic, non-commercial research  
184 purposes and in compliance with the Google Maps/Google Earth Additional Terms of Service. The released  
185 CrowdSat dataset does not include any Google Earth images. However, Google Earth imagery for some  
186 regions suffers from cloud or haze contamination, mosaic seams, compression artifacts or relatively sparse  
187 visible crowds, which limits its usefulness for dense crowd annotation. To complement these cases, we  
188 additionally used imagery from the BJ3N and JL04A satellites. The BJ3N satellite operates in a  
189 Sun-synchronous orbit at an altitude of 610 km, with a 5-day revisit period and an LTAN of 11:00. It captures  
190 RGB and panchromatic bands with a spatial resolution of 1.20 m and 0.30 m, respectively. The JL04A satellite  
191 follows a Sun-synchronous orbit at 535 km, with a shorter 3-day revisit period and an LTAN at 10:30,  
192 allowing for more frequent observations. The JL04A imagery used in CrowdSat has a spatial resolution of 0.31  
193 m for the panchromatic band and 1.24 m for the multispectral bands.

194



195

196 Fig. 3. Spatial distribution map of collected VFR satellite imagery for some locations where crowds typically gather. Part of the  
 197 imagery displayed in this figure was obtained from the Google Earth platform (© Google Earth 2025).

198

199 To ensure broad spatial coverage and diverse conditions, we selected imagery spanning 32 provincial-level  
 200 divisions in China (except Guizhou Province and Macao, due to fewer satellite images in these areas). Our  
 201 dataset collection follows two main principles. First, the collected imagery must be visually reliable for  
 202 annotation, i.e., cloud- or haze-free, with minimal shadow contamination and high radiometric quality. Second,  
 203 the selected regions need to contain a sufficient number of clearly visible individuals. Within each  
 204 provincial-level division, we first identified typical public gathering places (e.g., city squares, commercial  
 205 streets, transport hubs and scenic sites), then reviewed imagery from BJ3N, JL04A satellites and the Google  
 206 Earth platform manually, and selected the scenes that best satisfy these two principles. These regions exhibit  
 207 substantial heterogeneity, encompassing various landscapes and urban settings, as illustrated in Fig. 3. These  
 208 scenes include open public areas (e.g., the Forbidden City), built-up areas (e.g., Kashgar People's Square),  
 209 snowy regions (e.g., Harbin Ice and Snow World), areas with lush vegetation (e.g., the Emperor Qinshihuang's  
 210 Mausoleum Site Museum), beaches (e.g., Qingdao Jiaozhou Bay National Marine Park), desert regions (e.g.,

211 Yuesha Island), etc.. All the imagery was acquired between Feb. 20, 2023, and Jan. 2, 2025. In addition, before  
212 using the collected imagery, preprocessing steps such as geometric correction, clipping and radiometric  
213 normalization were performed to enhance the consistency between images from different sources.

214 Using the original 1.20 m BJ3N or 1.24 m JL4A multi-spectral imagery directly makes it challenging to  
215 distinguish individuals. To address this limitation, a state-of-the-art pan-sharpening technique known as  
216 area-to-point regression kriging (ATPRK) (Wang et al., 2015, 2016), which attained smaller relative  
217 global-dimensional synthesis error (ERGAS) and larger spectral angle mapper (SAM) values than classical  
218 fusion methods, was employed to increase the spatial resolution of the BJ3N and JL4A imagery to 0.30 m and  
219 0.31 m, respectively. The fused imagery with finer spatial resolution enables a more accurate interpretation for  
220 subsequent analysis.

221

## 222 2.2. Data Labeling

223

224 Labeling individuals in VFR satellite imagery presents unique challenges compared to ground and aerial  
225 imagery. Due to the relatively blurred edges of individuals in VFR satellite imagery, it is difficult to delineate  
226 their contours. This increases the likelihood of mistaking other ground objects as individuals. For example, as  
227 shown in Fig. 4(a), objects within the red and cyan circles exhibit visual characteristics similar to individuals.  
228 However, the objects in the red circles correspond to road asphalt, while those in the cyan circles are street  
229 lamps. Similarly, other stationary objects, such as stone pillars, may also be labeled incorrectly, which can  
230 further mislead network training.

231



Fig. 4. Examples of mislabeling. (a) and (b) were captured in the same region on Feb. 16, 2023 from the BJ3N satellite and Aug. 7, 2024 from the Google Earth platform (© Google Earth 2025), respectively. The objects in the red and cyan circles correspond to road asphalt and street lamps, respectively.

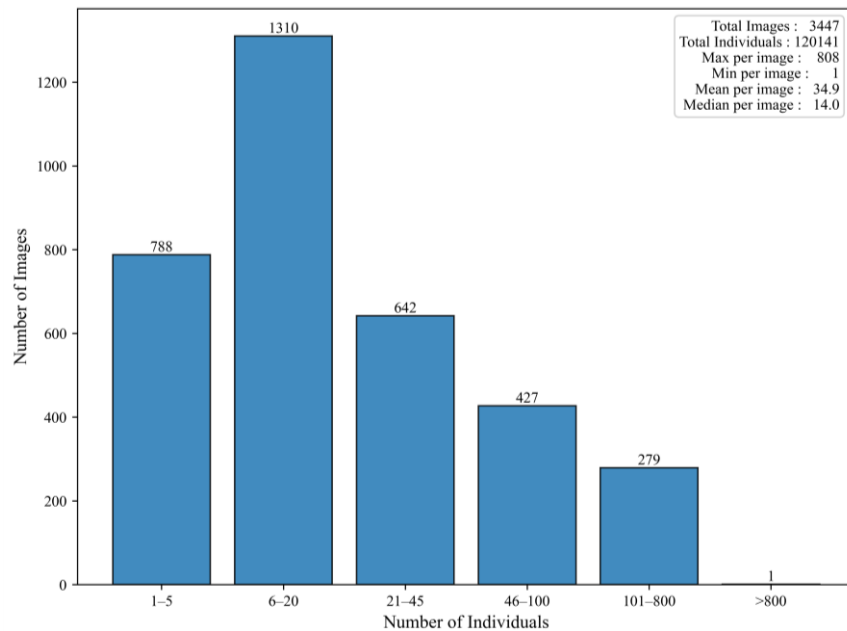
To alleviate the above issue, multi-temporal VFR imagery from the Google Earth platform was used as auxiliary data, due to the extensive spatial and temporal coverage of Google Earth imagery (Li et al., 2023; Shi et al., 2023; Ni et al., 2024). For example, Fig. 4(b) was captured in the same region at a different time, making it easier to distinguish individuals from other small objects. By cross-referencing imagery captured at different time points, fixed objects were identified as non-individuals, reducing labeling errors. Moreover, to further ensure accurate annotations, a center-point labeling strategy was employed. Specifically, for each individual, a single point was placed at the center of the  $3 \times 3$  pixel region to present the approximate position in the imagery. This approach prevents boundary ambiguity and minimizes the risk of mislabeling non-human objects.

To facilitate the use of the data in a deep learning architecture, following the collection of classical remote sensing datasets, such as LEVIR-CD (Chen and Shi, 2020) and WHU-CD (Ji et al., 2018), the collected imagery was cropped into  $256 \times 256$  pixel patches without overlapping. Patches without individuals were removed, resulting in a total of 3,447 labeled patches.

### 2.3. Data Analysis

252 Compared to ground and aerial imagery, CrowdSat sets new benchmarks with its vast scale, diverse density  
 253 distribution and comprehensive multi-environment coverage. It presents four key characteristics.

254 *1) Multi-Density Representation.* As illustrated in Fig. 5, the CrowdSat dataset consists of a total of 3,447  
 255 patches, covering 120,141 individuals. The number of individuals in each patch varies significantly, with a  
 256 maximum of 808 individuals and a minimum of 1 individual. On average, each patch contains 34.9 individuals,  
 257 while the median number per patch is 14.0. This wide range of crowd densities across different patches offers  
 258 a comprehensive representation of various crowd scenarios, making the dataset valuable for evaluating crowd  
 259 detection methods under diverse conditions.

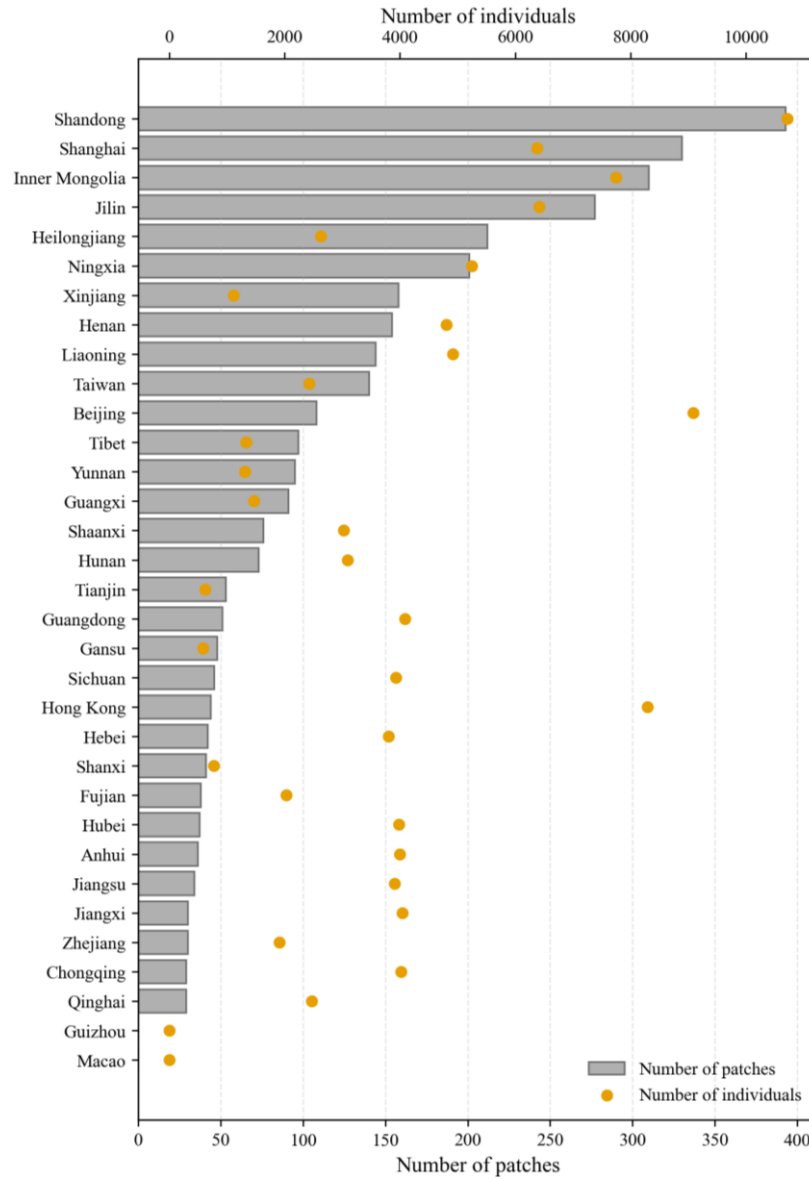


261  
 262 Fig. 5. Crowd count distribution in CrowdSat.

263  
 264 *2) Extensive National Coverage.* CrowdSat spans China (except Guizhou Province and Macao) and  
 265 includes samples from 32 Chinese provincial-level divisions. As illustrated in Fig. 6, the number of patches  
 266 per province ranges from 0 to 393, and the corresponding number of individuals ranges from 0 to 10,716.  
 267 Provinces with frequent large public gatherings, such as Shandong (393 patches; 10,716 individuals) and  
 268 Shanghai (330 patches; 6,374 individuals), contribute more samples, whereas sparsely populated or remote

269 regions, such as Tibet and Qinghai, contain fewer but still representative scenes. Guizhou and Macao currently  
270 have no labelled patches that satisfy our selection principles (cloud- and haze-free, minimal shadow  
271 contamination and a sufficient number of visible individuals). The vast coverage of CrowdSat provides a wide  
272 range of crowd types, promoting the generalizability of the dataset to various scenarios.

273



274

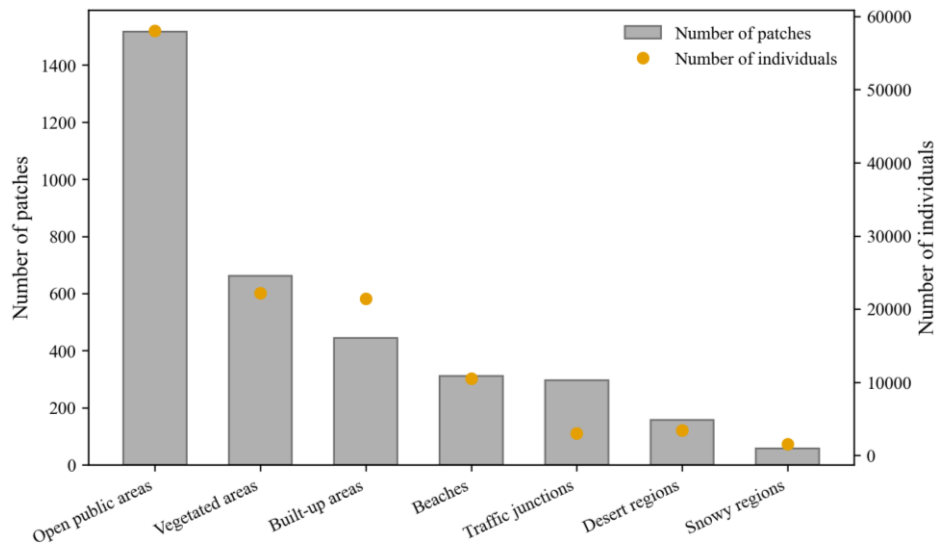
275 Fig. 6. Distribution of patches and individuals per province in CrowdSat.

276

277 3) *Multi-Season Adaptation*. CrowdSat was collected between Feb. 20, 2023 and Jan. 2, 2025, covering all  
 278 four seasons. This wide temporal span allows models to adapt inherently to seasonal variations in terms of  
 279 lighting conditions (e.g., summer glare and winter haze) and crowd movement patterns.

280 4) *Diverse Environmental Representations*. CrowdSat contains samples from distinct environmental  
 281 categories across China, such as open areas, snowy regions, beaches, desert regions, etc.. As shown in Fig. 7,  
 282 CrowdSat is dominated by open public areas, where mass gatherings occur commonly. Built-up areas and  
 283 vegetated areas together account for most of the remaining samples, consisting of dense commercial streets,  
 284 residential blocks and urban parks where crowd activities take place everyday. Beaches form a medium-sized  
 285 component of the dataset, whereas snowy regions and desert regions are less frequent but still represented by  
 286 dozens of patches and several thousand individuals. Traffic junctions contain a noticeable number of patches  
 287 but fewer individuals per patch. By systematically encompassing scenarios from historical landmarks to  
 288 modern transportation hubs, this cross-environment integration can alleviate location-specific bias.

289



290

291 Fig. 7. Distribution of patches and individuals across different scene types in CrowdSat.

292

293

### 294 3. Methods

295

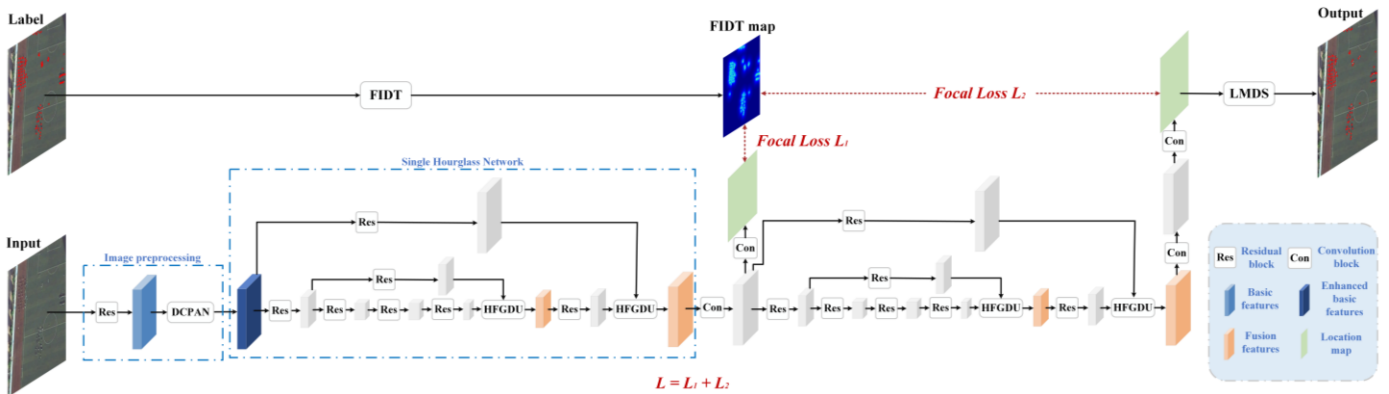
296 *3.1. Overview of CrowdSat-Net*

297

298 To detect crowds in large-scale VFR satellite imagery effectively, a novel point-based CNN, CrowdSat-Net,  
 299 was proposed. It was specifically designed to address two key challenges: 1) gradual blurring and loss of small  
 300 object signals during feature extraction and 2) high-frequency information loss during upsampling operations.

301 An overview of the CrowdSat-Net architecture is illustrated in Fig. 8. At the beginning, the label image is  
 302 transformed into a Focal Inverse Distance Transform (FIDT) map (Liang et al., 2022a) using the FIDT method.  
 303 CrowdSat-Net adopts a two-stage stacked Hourglass Network (Newell et al., 2016) (a typical network used for  
 304 point feature extraction), with key modifications to increase the accuracy of CD. First, to improve small-object  
 305 feature representation, a novel module, that is, DCPAN, is embedded in the image preprocessing stage. Second,  
 306 during the upsampling operation, the traditional method (e.g., bilinear interpolation) is replaced with a novel  
 307 module, HFGDU, which aims to restore lost high-frequency information. Within each Hourglass Network, a  
 308 location map is generated, and the loss is computed using the FIDT map and the location map through the  
 309 Focal Loss function (Lin et al., 2017b), a strategy known as Intermediate Supervision. The total loss  $L$  is  
 310 computed as the sum of Focal Losses from each Hourglass Network. Finally, the location map from the last  
 311 Hourglass Network is processed using a Local-Maxima-Detection-Strategy (LMDS) to obtain the final  
 312 individual localizations.

313



314

315 Fig. 8. Overview of the proposed CrowdSat-Net. First, the labeled image is transformed into the FIDT map. During each training  
 316 iteration, CrowdSat-Net enhances the basic features in the image preprocessing using the DCPAN module. Then, these enhanced  
 317 features pass through the two-stacked Hourglass Network. In each Hourglass Network, the traditional upsampling method is  
 318 replaced with the HFGDU module to recover the lost fine details. Each Hourglass Network generates a location map, which is  
 319 compared with the FIDT map to calculate the Focal Loss. The total loss  $L$  is the sum of Focal Losses from each Hourglass Network.  
 320 After training, the location map conducted by the last Hourglass Network is transformed into the actual localization result using the  
 321 LMDS method.

322

323 The remainder of this section introduces important components in detail, including the FIDT map, DCPAN,  
 324 HFGDU, LMDS and model evaluation.

325

### 326 3.2. Focal Inverse Distance Transform (FIDT) Map

327

328 It is critical to identify individual localization in CrowdSat-Net. Traditional methods, such as binary-like  
 329 maps (Liu et al., 2019), segmentation-like maps (Xu et al., 2022), topological maps (Abousamra et al., 2021)  
 330 and independent instance maps (Gao et al., 2020c), struggle to distinguish overlapping objects in dense crowds  
 331 due to their reliance on fixed thresholds or semantic boundaries. To overcome this problem, the FIDT (Liang  
 332 et al., 2022a) method was employed, which enables overlap-free head localization by assigning higher pixel  
 333 responses closer to head centers, thereby ensuring accurate localization in dense crowds. The FIDT is defined  
 334 as follows:

335

$$336 \quad I = \frac{1}{d(x,y)^{(\alpha \times P(x,y) + \beta) + C}} \quad (1)$$

337

338 where  $I$  represents the FIDT map,  $\alpha$  and  $\beta$  are weight factors, set to the default values of 0.02 and 0.75,  
 339 respectively (Liang et al., 2022a).  $C$  indicates a constant, which aims to prevent division by zero.  $d(x,y)$

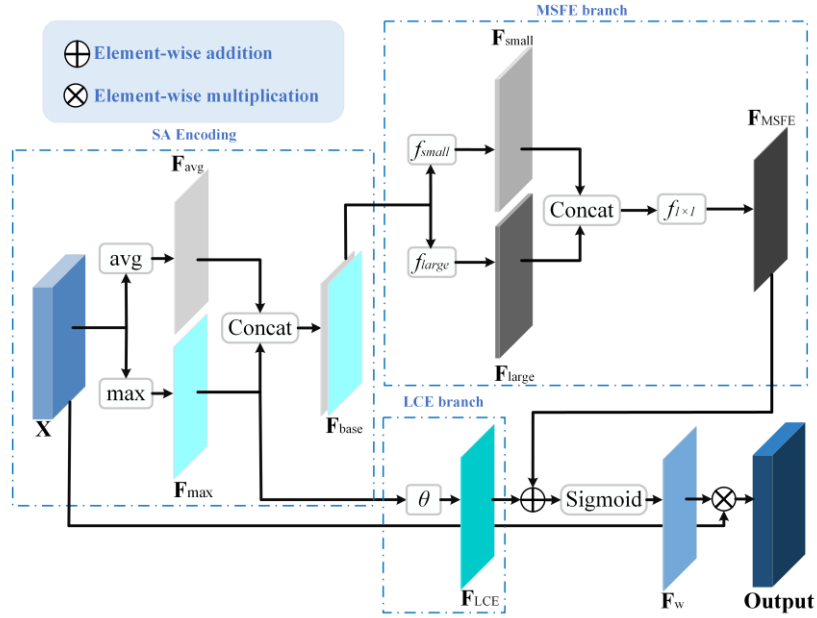
denotes the distance between the pixel at location  $(x, y)$  and its nearest annotated pixel at location  $(x', y')$ , in the set of all annotation  $\mathbf{B}$ :

$$d(x, y) = \min_{(x', y') \in \mathbf{B}} \sqrt{(x - x')^2 + (y - y')^2} \quad (2)$$

### 3.3. Dual-Context Progressive Attention Network (DCPAN)

The image preprocessing stage aims to extract fundamental feature representations from the original inputs, serving as the foundation for downstream network architectures. Conventional implementations in the Hourglass Network employ typically a standard convolutional layer followed by a ResNet layer (Zhou et al., 2019), which ignores the gradual signal attenuation of small objects. This limitation can compromise model sensitivity to fine-grained features and increase false negative rates in individual detection.

To address the above issue, we first excluded the standard convolutional layer and then proposed the DCPAN module embedded after the ResNet layer, shown in Fig. 9. The DCPAN module combines synergistically a base spatial attention (SA) Encoding and two parallel branches (multi-scale feature extraction (MSFE) and local contrast enhancement (LCE) branches), aiming to improve the small-object signal presentation. Specifically, given a feature map  $\mathbf{X} \in \mathbb{R}^{C \times H \times W}$  generated by the ResNet layer, where  $C$ ,  $H$  and  $W$  denote the channel, height and width of the feature, respectively, the DCPAN module generates small-object feature enhancement weights through three key operations detailed as follows.



360

361 Fig. 9. Flowchart of the DCPAN module, which consists of three main components: SA Encoding, MSFE and LCE branches. The  
 362 SA Encoding component extracts the base SA feature from the original features  $\mathbf{X}$ . The MSFE branch captures contextual  
 363 information from this base feature, while the LCE branch focuses on extracting its local contrast. Finally, the outputs of the MSFE  
 364 and LCE branches are fused through element-wise addition, and the result is passed through a sigmoid activation function to obtain  
 365 the enhancement weight map  $\mathbf{F}_w$ . This map is subsequently applied to the original feature  $\mathbf{X}$  via element-wise multiplication.

366

367 1) *SA Encoding*. Small objects lack distinct semantic features, but exhibit unique spatial patterns. We first,  
 368 thus, extracted the base spatial attention cues  $\mathbf{F}_{base}$ , achieved as follows:

369

$$370 \quad \mathbf{F}_{max} = \max(\mathbf{X}), \mathbf{F}_{avg} = \frac{1}{C} \sum_{c=1}^C \mathbf{X}_c \quad (3)$$

371

372 where  $\mathbf{F}_{max} \in \mathbb{R}^{1 \times H \times W}$  and  $\mathbf{F}_{avg} \in \mathbb{R}^{1 \times H \times W}$  represent the max-pooled and average-pooled features,  
 373 respectively. These two features are then concatenated to form the base spatial feature  $\mathbf{F}_{base} \in \mathbb{R}^{2 \times H \times W}$ .

374 2) *MSFE Branch*. Although the base spatial feature can improve the presentation of prominent objects by  
 375 highlighting global significant regions, the global pooling operations (e.g., max/avg) in SA Encoding fail to  
 376 model positional dependencies between objects, leading to suboptimal results in cluttered scenes (Fu et al.,

377 2019). To alleviate this issue, a simple MSFE branch, which employs parallel dilated convolutional layers with  
 378 complementary dilation rates, was proposed:

379

$$380 \quad \mathbf{F}_{\text{small}} = f_{\text{small}}(\mathbf{F}_{\text{base}}), \mathbf{F}_{\text{large}} = f_{\text{large}}(\mathbf{F}_{\text{base}}) \quad (4)$$

381

382 where  $f_{\text{small}}$  and  $f_{\text{large}}$  are the dilated convolutional operations with dilation rates of 2 and 4, respectively.

383 Then, the  $\mathbf{F}_{\text{MSFE}} \in \mathbb{R}^{1 \times H \times W}$  of these layers are concatenated and fused using a  $1 \times 1$  convolutional layer:

384

$$385 \quad \mathbf{F}_{\text{MSFE}} = f_{1 \times 1}(\text{Concat}(\mathbf{F}_{\text{small}}, \mathbf{F}_{\text{large}})) \quad (5)$$

386

387 3) *LCE Branch*. Standard pooling operations always blur high-frequency details, leading to degraded  
 388 accuracy on fine-grained detection tasks (Fu et al., 2019). To address this problem, the LCE branch highlights  
 389 locations with large heterogeneity by a local contrast generator  $\theta$ :

390

$$391 \quad \mathbf{F}_{\text{LCE}} = \theta(\mathbf{F}_{\text{max}}) = \partial(|\mathbf{F}_{\text{max}} - \text{AvgPool}_{3 \times 3}(\mathbf{F}_{\text{max}})|) \quad (6)$$

392

393 where  $\mathbf{F}_{\text{LCE}} \in \mathbb{R}^{1 \times H \times W}$  represents the contrast-sensitive weight map and  $|\cdot|$  denotes absolute value  
 394 computation.  $\partial$  is the learnable contrast enhancement operator composed of a  $3 \times 3$  convolutional layer, batch  
 395 normalization, ReLU activation and a  $3 \times 3$  convolutional layer.

396 Finally, the enhancement weight map  $\mathbf{F}_w$  is generated by applying a sigmoid activation to the element-wise  
 397 addition of  $\mathbf{F}_{\text{MSFE}}$  and  $\mathbf{F}_{\text{LCE}}$ . This weight map is then multiplied pixel-wise with the original feature  $\mathbf{X}$  to obtain  
 398 the final feature.

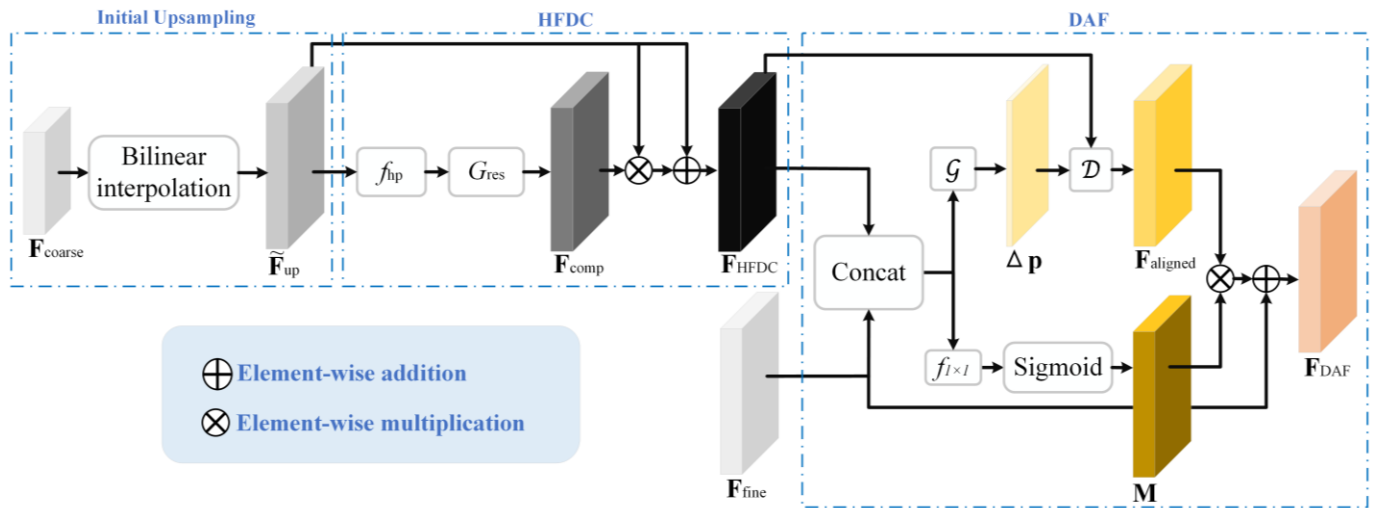
399

#### 400 3.4. High-Frequency Guided Deformable Upsampler (HFGDU)

401

402 Within each Hourglass network, the coarse spatial resolution features are upsampled and then fused with  
 403 fine spatial resolution features to increase the ability to detect small objects. However, as mentioned in the  
 404 Introduction, traditional upsampling methods, such as nearest neighbor and bilinear interpolation, often fail to  
 405 recover fine spatial details, leading to over-smooth boundaries (Li et al., 2024) and misalignment of  
 406 high-frequency features (Dai et al., 2017), making it challenging to detect small objects effectively.

407 To alleviate the above challenge, the HFGDU module was introduced. As shown in Fig. 10, HFGDU  
 408 alleviates the dual challenges of high-frequency detail recovery and geometric-aware feature alignment  
 409 through a three-stage architecture: Initial Upsampling, High-Frequency Detail Compensation (HFDC) and  
 410 Deformable Alignment Fusion (DAF).



412  
 413 Fig. 10. The structure of the HFGDU module. It contains three main components: Initial Upsampling, HFDC and DAF. Initial  
 414 Upsampling uses the bilinear interpolation to match the spatial resolution of coarse and fine features. HFDC recovers the  
 415 high-frequency details of the upsampled features, while DAF refines the spatial alignment of the high-frequency features.

416  
 417 1) *Initial Upsampling*. Given an input coarse spatial resolution feature map  $\mathbf{F}_{\text{coarse}} \in \mathbb{R}^{C' \times H' \times W'}$  and a fine  
 418 spatial resolution feature map  $\mathbf{F}_{\text{fine}} \in \mathbb{R}^{C' \times 2H' \times 2W'}$ , to match the spatial size, we first applied the simple  
 419 bilinear interpolation to  $\mathbf{F}_{\text{coarse}}$  to obtain an upsampled feature map  $\tilde{\mathbf{F}}_{\text{up}} \in \mathbb{R}^{C' \times 2H' \times 2W'}$ . However, bilinear  
 420 interpolation always smooths feature boundaries, failing to recover the lost high-frequency details that are

421 critical for small object detection.

422 2) *HFDC*. To recover the lost fine details during upsampling, we introduced a learnable high-pass filter  $f_{\text{hp}}$   
423 to extract high-frequency components from the upsampled feature map:

424

$$425 \quad \mathbf{F}_{\text{hf}} = f_{\text{hp}}(\tilde{\mathbf{F}}_{\text{up}}) \quad (7)$$

426

427 where the filter  $f_{\text{hp}}$  is designed based on a Laplacian-like kernel, which is initialized with a center weight of 1  
428 and surrounding weights of  $-1/8$ , to enhance edge structures. The extracted high-frequency feature map  $\mathbf{F}_{\text{hf}}$  is  
429 then refined using a residual compensation generator  $G_{\text{res}}$ :

430

$$431 \quad \mathbf{F}_{\text{comp}} = G_{\text{res}}(\mathbf{F}_{\text{hf}}) \quad (8)$$

432

433 where  $G_{\text{res}}$  is composed of a  $3 \times 3$  convolutional layer followed by two additional  $3 \times 3$  convolutional layers  
434 with ReLU activations, and a final  $3 \times 3$  convolutional layer with a Sigmoid activation to constrain the  
435 compensation magnitude. The map  $\mathbf{F}_{\text{comp}}$  is then integrated with the upsampled feature map to generate the  
436 compensated high-frequency details map  $\mathbf{F}_{\text{HFDC}}$ :

437

$$438 \quad \mathbf{F}_{\text{HFDC}} = \tilde{\mathbf{F}}_{\text{up}} \cdot (\mathbf{1} + \mathbf{F}_{\text{comp}}) \quad (9)$$

439

440 where  $\cdot$  means element-wise multiplication. This step helps in recovering the fine high-frequency components,  
441 which can increase potentially the clarity and sharpness of the upsampled feature representation.

442 3) *DAF*. After compensating for high-frequency details, we employed DAF to refine the spatial alignment  
443 between the upsampled feature  $\mathbf{F}_{\text{HFDC}}$  and the fine spatial resolution feature  $\mathbf{F}_{\text{fine}}$ . Traditional upsampling  
444 methods rely on fixed-grid sampling, which may cause spatial misalignment. To alleviate this, we introduced a  
445 deformable convolution layer  $\mathcal{D}$  that predicts spatial offsets  $\Delta \mathbf{p}$ :

446

447

$$\mathbf{F}_{\text{aligned}} = \mathcal{D}(\mathbf{F}_{\text{HFDC}}, \Delta\mathbf{p}) \quad (10)$$

448

449 where  $\mathbf{F}_{\text{aligned}}$  denotes the deformably sampled feature map used to correct misalignment, and the offsets  $\Delta\mathbf{p}$   
 450 are adaptively learned from the feature map:

451

452

$$\Delta\mathbf{p} = \mathcal{G}(\text{Concat}(\mathbf{F}_{\text{HFDC}}, \mathbf{F}_{\text{fine}})) \quad (11)$$

453

454 In Eq. (11),  $\mathcal{G}$  is a lightweight offset prediction network composed of two stacked  $3 \times 3$  convolutional layers  
 455 that estimates local displacements. This helps the model to align features dynamically based on the underlying  
 456 structure.

457

458 Furthermore, to integrate selectively aligned features while suppressing redundant information, we  
 459 introduced a feature modulation gate  $M$  that acts as an adaptive attention mechanism:

459

460

$$\mathbf{F}_{\text{DAF}} = \mathbf{M} \odot \mathbf{F}_{\text{aligned}} + \mathbf{F}_{\text{fine}} \quad (12)$$

461

462 where  $\mathbf{F}_{\text{DAF}}$  is the final fusion map and  $\odot$  denotes element-wise multiplication, and the modulation gate  $\mathbf{M}$  is  
 463 computed as:

464

465

$$\mathbf{M} = \sigma(f_{1 \times 1}(\text{Concat}(\mathbf{F}_{\text{HFDC}}, \mathbf{F}_{\text{fine}}))) \quad (13)$$

466

467 where  $\sigma$  is the sigmoid activation. By integrating these three components, HFGDU aims to preserve  
 468 high-frequency details and dynamically align features effectively, which can contribute to improving the  
 469 detection of small objects in the Hourglass Network.

470

### 3.5. Local Maxima Detection Strategy

The FIDT map generated by the last Hourglass Network can identify the center point of each individual effectively, but it struggles with filtering false positives. To alleviate this issue, the LMDS was introduced to determine the correct position of each individual (Liang et al., 2022a). Specifically, a  $3 \times 3$  max-pooling operation is employed to generate all candidate points, and an adaptive threshold, denoted as  $\delta$ , is used to filter out false positives. This threshold is defined empirically as  $100 / 255.0$  times the maximum value of the FIDT map, and only points with values greater than or equal to  $\delta$  are selected. Additionally, if the maximum value of the FIDT map is smaller than a tiny fixed threshold (set to 0.10) (Liang et al., 2022a), it indicates that there are no detected individuals in the input image.

### 3.6. Model Evaluation

To analyze the performance of the proposed CrowdSat-Net model comprehensively, both localization and counting metrics are introduced in this research. For localization metrics, the first step in assessing localization performance is to match the predicted individual point  $\hat{P}$  with the ground reference individual point  $P$ . If the point-marching distance between  $\hat{P}$  and  $P$  is less than a threshold  $\gamma$ , they are considered matched. In this study,  $\gamma$  was set to 2 pixels, and the nearest neighbors method (Kramer, 2013) was employed for matching, with each point being assigned to its single closest neighbor. After matching, the match matrix and counts of the number of True Positive (TP), False Positive (FP) and False Negative (FN) are obtained. Then, Precision, Recall, F1-score (Wang et al., 2020a) and Mean Localization Error (MLE) (Sam et al., 2021) were adopted based on TP, FP, FN and point-marching distance to evaluate the performance of the proposed model. The formulae are as follows:

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (14)$$

496 
$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (15)$$

497 
$$\text{F1-score} = \frac{2 \times \text{Recall} \times \text{Precision}}{\text{Recall} + \text{Precision}} \quad (16)$$

498 
$$\text{MLE} = \frac{1}{N_{\text{TP}}} \sum_{j=1}^{N_{\text{TP}}} \|\hat{\mathbf{P}}_j - \mathbf{P}_j\|_2 \quad (17)$$

499

500 where  $\{\hat{\mathbf{P}}_j\}_{j=1, \dots, N_{\text{TP}}}$  and  $\{\mathbf{P}_j\}_{j=1, \dots, N_{\text{TP}}}$  denote the coordinates of the matched predicted and ground-reference  
 501 points, respectively, and  $N_{\text{TP}}$  represent the number of TP.  $\|\hat{\mathbf{P}}_j - \mathbf{P}_j\|_2$  corresponds to the Euclidean distance  
 502 between the predicted and ground-reference points. For counting metrics, the Mean Absolute Error (MAE)  
 503 and Root Mean Squared Error (RMSE) were included, which are defined as follows:

504

505 
$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |C_i^{\text{pred}} - C_i^{\text{re}}| \quad (18)$$

506 
$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (C_i^{\text{pred}} - C_i^{\text{re}})^2} \quad (19)$$

507

508 where  $C_i^{\text{pred}}$  and  $C_i^{\text{re}}$  are the predicted and ground-reference counts in the  $i$ -th image patch, respectively.

509

510

## 511 4. Experiments

512

513 In this section, comprehensive experiments on the Crowd-Sat dataset were conducted to examine the  
 514 effectiveness of the proposed CrowdSat-Net method. Specifically, Section 4.1 details the experimental setup.  
 515 Section 4.2 presents a series of ablation studies to validate the contributions of the DCPAN and HFGDU  
 516 modules. Subsequently, Section 4.3 compares the CrowdSat-Net method with eight recent state-of-the-art CD  
 517 benchmark methods (originally designed for ground or aerial imagery and satellite-based animal detection) to  
 518 assess its performance. Section 4.4 examines the sensitivity of localization performance to the point-matching

519 threshold. Section 4.5 presents the computational efficiency and memory usage of different CD models.  
520 Section 4.6 illustrates the performance of various backbones in CrowdSat-Net. Section 4.7 analyzes the  
521 performance of CrowdSat-Net across different crowd density levels. Finally, Section 4.8 examines the  
522 cross-regional generalization of the CrowdSat-Net method by applying it to unseen foreign regions.

#### 524 4.1. Experimental Setup

525  
526 In this research, the labeled patches in CrowdSat were divided randomly into training and validation sets in  
527 a 4:1 ratio. To increase the robustness and generalization of models, we applied augmentation techniques,  
528 including horizontal flipping, vertical flipping and CutMix (Yun et al., 2019). After augmentation, patches  
529 without individuals were removed, resulting in 10,357 training patches and 704 validation patches.

530 After preparing the dataset, we trained the CrowdSat-Net model, all baseline models and other CD models  
531 for 150 epochs using the Adam optimizer (Adam et al., 2014) with a base learning rate of 0.0003 and a weight  
532 decay of 0.001, on the same augmented dataset described above to ensure fair comparison. For all models, the  
533 batch size was set to 8, and they were implemented in PyTorch 1.21. All experiments were conducted on a  
534 Windows 10 workstation equipped with a 13th Gen Core i7-13700 CPU and a 24 GB NVIDIA GeForce RTX  
535 4090 GPU.

#### 536 537 4.2. Ablation Experiments

##### 538 539 4.2.1. Ablation of DCPAN and HFGDU.

540  
541 To examine the contributions of the DCPAN and HFGDU modules in CrowdSat-Net, a basic ablation  
542 experiment was conducted. In this experiment, the baseline refers to CrowdSat-Net with the DCPAN module  
543 removed and the HFGDU module replaced by bilinear interpolation.

544 As presented in Table 1, introducing either DCPAN or HFGDU outperforms the baseline in terms of  
 545 localization and counting metrics. For example, with DCPAN alone, the F1-score and Precision increase by  
 546 0.90% and 4.36%, respectively, while the MLE decreases by 0.02. The MAE and RMSE decrease by 0.83 and  
 547 2.74, respectively. A similar trend is observed when only HFGDU is enabled, that is, the F1-score rises to  
 548 65.52%, and MLE, MAE and RMSE are reduced to 0.40, 12.69 and 35.55, respectively. Notably, combining  
 549 both modules achieves the best overall performance, with F1-score and Precision increased, and MLE, MAE  
 550 and RMSE decreased compared with the baseline.

551

552

Table 1 Basic ablation study for DCPAN and HFGDU.

Methods	DCPAN	HFGDU	F1-score (%)	Recall (%)	Precision (%)	MLE	MAE	RMSE
Baseline	×	×	64.42	60.77	68.54	0.42	13.04	37.94
	√	×	65.32	59.16	72.90	0.40	12.21	35.20
CrowdSat-Net	×	√	65.52	<b>61.10</b>	70.63	0.40	12.69	35.55
	√	√	<b>66.12</b>	60.27	<b>73.23</b>	<b>0.39</b>	12.07	34.37

553

554

#### 4.2.2. Ablation Study of Different Feature Enhancement Modules in CrowdSat-Net.

555

556

557

558

559

560

561

562

563

To further examine the effectiveness of DCPAN, a comparison with two well-known feature enhancement modules (spatial attention (SA) (Wang et al., 2017) and efficient channel attention (ECA) (Wang et al., 2020b)) was performed. In this experiment, all components were the same as the baseline method except for the feature enhancement module. The result is displayed in Table 2. It is seen that DCPAN achieves greater performance than SA and ECA, achieving the largest F1-score of 65.32% (+0.46% over SA, +0.58% over ECA) and Recall of 59.16% (+1.32% over SA, +1.02% over ECA). It also produces the smallest MAE (12.21) and RMSE (35.20), decreasing MAE by 0.93 and 1.23 and RMSE by 2.67 and 3.49 compared to SA and ECA, respectively.

564 Architecturally, DCPAN extends SA by integrating two novel branches: the MSFE and LCE. Specifically,  
 565 MSFE utilizes dilated convolutions for multi-scale contextual information, mitigating false negatives  
 566 effectively caused by the fixed receptive field of SA. Furthermore, in contrast to SA, the LCE enhances feature  
 567 representation by amplifying local feature responses through contrast-based spatial attention. These  
 568 innovations increase the localization accuracy of DCPAN collectively. Compared to ECA, which focuses on  
 569 channel-wise attention exclusively, DCPAN alleviates spatial signal blurring through its dual-branch design.  
 570

571 Table 2 Comparison of DCPAN with SA and ECA in CrowdSat-Net.

Methods	F1-score (%)	Recall (%)	Precision (%)	MLE	MAE	RMSE
SA (Wang et al., 2017)	64.86	57.84	<b>73.82</b>	<b>0.38</b>	13.14	37.87
ECA (Wang et al., 2020b)	64.74	58.14	73.03	0.41	13.44	38.69
DCPAN	<b>65.32</b>	<b>59.16</b>	72.90	0.40	<b>12.21</b>	<b>35.20</b>

572

### 573 4.2.3. Ablation of Different Upsamplers in CrowdSat-Net.

574

575 To reveal comprehensively the effectiveness of HFGDU, a comparative analysis with four established  
 576 benchmark upsamplers, including bilinear interpolation, Content-Aware ReAssembly of FEatures (CARAFE)  
 577 (Wang et al., 2019), Deconvolution (Long et al., 2015), Pixel-shuffle (Shi et al., 2016), Dense Upsampling  
 578 Convolution (DUC) (Zhao et al., 2017) and Dysample (Liu et al., 2023b), was conducted.

579 The result in Table 3 demonstrates that HFGDU exhibits greater localization performance over existing  
 580 upsamplers, achieving the largest F1-score and Precision. This is attributed to its dual-domain enhancement  
 581 mechanism: 1) The HFDC module employs a learnable Laplacian operator to amplify boundary gradients,  
 582 mitigating the inherent edge blurring in bilinear interpolation effectively and yielding a 0.33% Precision gain,  
 583 a 1.11% Recall increase and a 0.80% F1-score increase over bilinear interpolation. 2) The DAF mechanism  
 584 based on dynamic migration prediction reduces the feature misalignment error by jointly optimizing the spatial  
 585 correspondence between coarse and fine spatial resolution features. While Dysample achieves the largest

586 Recall, its Precision lags obviously behind HFGDU (-2.79%). While Deconvolution achieves a larger Recall,  
 587 its fixed transposed convolution kernels lead to checkerboard artifacts, decreasing Precision by 2.44%.  
 588 Although CARAFE generates upsampled cores via content-sensing, it struggles with small objects in dense  
 589 scenes, leading to over-smooth results and a lower F1-score (65.17%) than bilinear interpolation (65.32%).  
 590 Pixel-shuffle performs upsampling by first expanding the channel dimension and then rearranging the  
 591 channels into a finer spatial grid. However, in this small-object detection task, its accuracy gain is limited.  
 592 DUC yields slightly larger F1-score and Recall than bilinear interpolation, but the overall accuracy gain  
 593 remains limited and still smaller than for HFGDU.

594

595

Table 3 Comparison between different upsamplers in CrowdSat-Net.

Methods	F1-score (%)	Recall (%)	Precision (%)	MLE	MAE	RMSE
Bilinear interpolation	65.32	59.16	72.90	0.40	12.21	35.20
CARAFE (Wang et al., 2019)	65.17	59.47	72.08	<b>0.38</b>	13.07	36.92
Deconvolution (Long et al., 2015)	65.38	60.74	70.79	0.41	12.76	35.84
Dysample (Liu et al., 2023b)	65.62	<b>61.42</b>	70.44	0.42	13.68	38.86
Pixel-shuffle (Shi et al., 2016)	64.98	59.44	71.66	0.42	13.35	37.14
DUC (Zhao et al., 2017)	65.46	60.12	71.84	0.41	12.88	36.14
HFGDU	<b>66.12</b>	60.27	<b>73.23</b>	0.39	<b>12.07</b>	<b>34.37</b>

596

### 597 4.3. Comparison Between CrowdSat-Net and Benchmark Methods

598

599 This section aims to compare the performance of the proposed CrowdSat-Net with six state-of-the-art CD  
 600 methods (originally designed based on ground or aerial imagery): Point to Point Network (P2PNet) (Song et  
 601 al., 2021), SCALNet (Wang et al., 2021), Crowd Localization TRansformer (CLTR) (Liang et al., 2022b),  
 602 Point quEry Transformer (PET) (Liu et al., 2023a), Focal Inverse Distance Transform Map for Crowd  
 603 Localization (FIDTMCL) (Liang et al., 2022a) and Auxiliary Point Guidance Crowd Counting (APGCC)

604 (Chen et al., 2024). Also, two satellite-based animal detection approaches, HerdNet (Delplanque et al., 2023)  
605 and U-Net-based ensemble model (Wu et al., 2023) (hereafter, UNE), were included.

606 As illustrated in Figs. 11 and A.1, five representative scenes, including traffic junctions, snowfields, dense  
607 urban regions, desert regions and other common impervious regions, were selected from 704 validation  
608 patches to demonstrate the detection performance across various scenarios. It is seen that SCALNet exhibits  
609 noticeable missed detections, particularly in cluttered junctions and dense urban areas where individual signals  
610 are easily confused with crosswalk markings, shadows and building edges within the yellow circles.  
611 FIDTMCL and APGCC alleviate this issue to some extent and detect more individuals in complex scenes.  
612 However, some of their predictions remain blurred around object boundaries (such as buildings, trees and  
613 tents), which leads to missed individuals at the periphery of gatherings. HerdNet captures several prominent  
614 clusters of individuals, yet the predictions are biased to high-density crowds and ignore frequently isolated  
615 individuals along roads or fences. P2PNet produces sparse predictions and often detects only a few scattered  
616 individuals in regions that contain clear gathering. The transformer-based methods, i.e., CLTR and PET,  
617 achieve limited visualization performance and, in densely populated areas, almost fail to detect crowds while  
618 producing anomalous outputs. UNE and CrowdSat-Net achieve greater detection performance with fewer  
619 missed or false detections, and, in the high- and low-density areas, their predictions are more accurate than  
620 those of the other benchmark methods.

621 Quantitative analysis on all 704 validation patches is shown in Table 4. It is clear that CrowdSat-Net  
622 achieves state-of-the-art localization performance with an F1-score of 66.12%, a Precision of 73.23% and an  
623 MLE of 0.39. Compared to the second most accurate method (UNE), CrowdSat-Net increases the F1-score by  
624 0.80% and Precision by 6.38% and decreases the MLE by 0.36. Notably, CrowdSat-Net outperforms  
625 SCALNet by 5.71% in terms of F1-score, 8.41% in terms of Recall and 0.88 in terms of MLE, highlighting the  
626 effectiveness in small-object detection. Although APGCC achieves a larger Recall and the smallest MAE and  
627 RMSE, its Precision is 7.53% smaller than that of CrowdSat-NET, and its MLE is 0.75 larger, revealing  
628 inherent limitations in point-based auxiliary supervision for the precise localization of small objects. HerdNet

629 obtains an MLE of 0.50, but its F1-score is only 60.18%. P2PNet and CLTR present relatively greater counting  
 630 performance, yet their localization performance is limited.

631

632

Table 4 Accuracy of nine CD methods on the CrowdSat dataset.

Methods	F1-score (%)	Recall (%)	Precision (%)	MLE	MAE	RMSE
SCALNET (Wang et al., 2021)	60.38	51.86	72.26	1.27	13.31	34.18
P2PNet (Song et al., 2021)	49.60	49.17	50.04	1.14	10.17	38.24
CLTR (Liang et al., 2022b)	13.50	13.27	13.73	1.24	9.26	22.06
PET (Liu et al., 2023a)	11.36	12.54	10.41	1.57	15.64	41.48
FIDTMCL (Liang et al., 2022a)	64.41	60.78	70.80	0.98	12.67	34.76
APGCC (Chen et al., 2024)	64.34	63.05	65.70	1.14	<b>8.31</b>	<b>17.32</b>
HerdNet (Delplanque et al., 2023)	60.18	57.31	63.35	0.50	14.49	38.93
UNE (Wu et al., 2023)	65.32	<b>63.87</b>	66.85	0.75	9.12	23.44
CrowdSat-Net	<b>66.12</b>	60.27	<b>73.23</b>	<b>0.39</b>	12.07	34.37

633





635

636

637

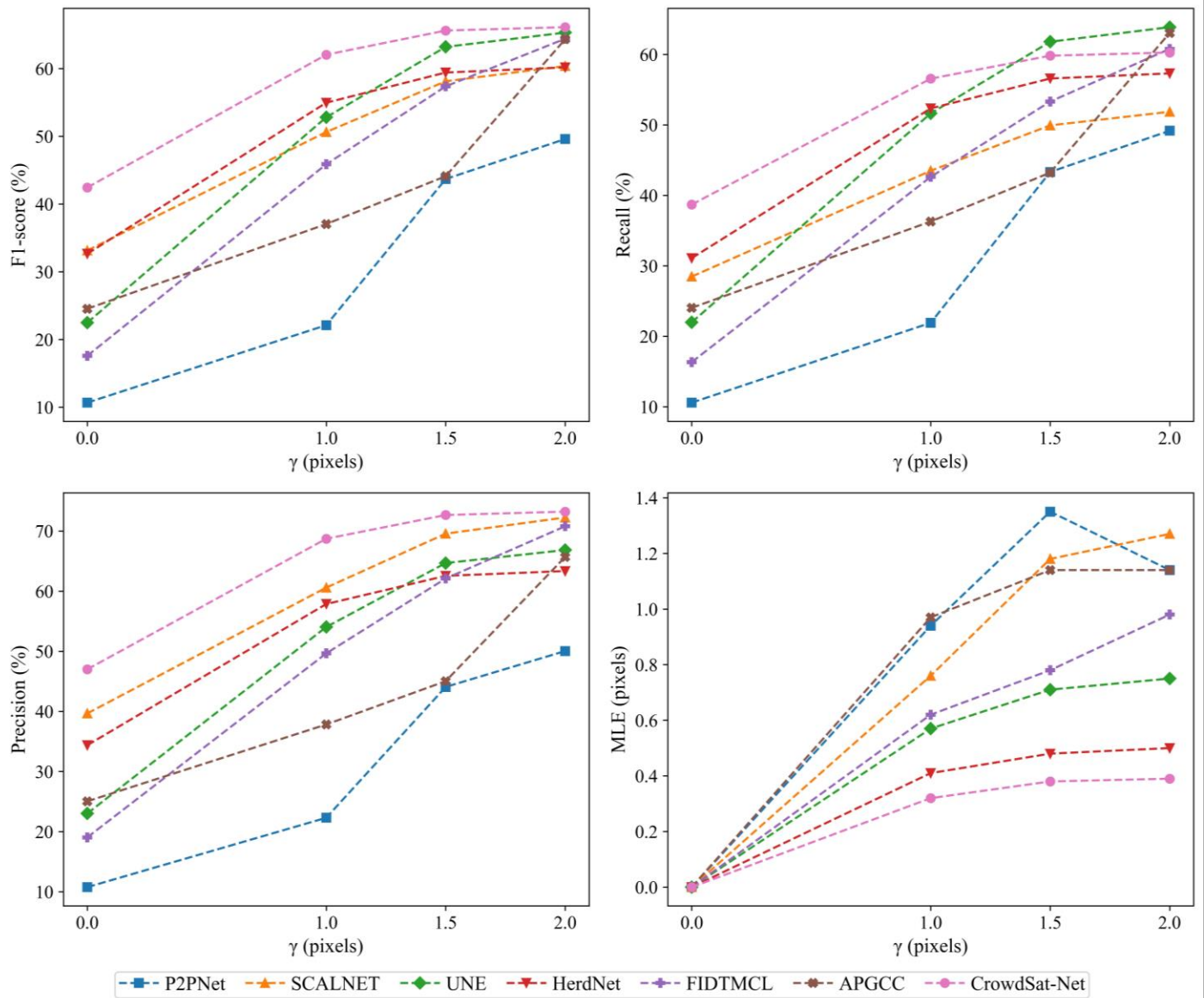
638

Fig. 11. Visual comparisons between different CD methods for the CrowdSat dataset. (a) Original image. (b) Reference. (c) SCALNet. (d) FIDTMCL. (e) APGCC. (f) HerdNet. (g) UNE. (h) CrowdSat-Net. (I) Traffic junctions. (II) Snowfields. (III) Dense urban regions. (IV) Desert regions. (V) Other common impervious regions. Some of the images in panels (I)-(V) were obtained from the Google Earth platform (© Google Earth 2025).

639

## 640 4.4. Sensitivity Analysis of Localization Performance to the Point-Matching Threshold

641



642

643 Fig. 12. Sensitivity of localization performance to the point-matching threshold across different CD methods.

644

645 To further examine the influence of the point-matching threshold, we conducted a sensitivity analysis with  
 646 respect to the threshold  $\gamma$ . All methods (except for CLTR and PET) were evaluated under the selected values,  
 647 i.e., 0.0, 1.0, 1.5 and 2.0. As illustrated in Fig. 12, CrowdSat-Net exhibits greater localization performance  
 648 over all benchmark methods in two aspects. First, CrowdSat-Net generally achieves the greatest accuracy

649 under all values of  $\gamma$ . Notably, when  $\gamma$  is set to 0, which means that a prediction is counted as correct only if its  
 650 pixel coordinates exactly coincide with those of a ground-reference point, CrowdSat-Net still obtains the  
 651 largest F1-score of 42.44%, Recall of 38.69% and Precision of 47.00%. These values exceed the second most  
 652 accurate results by 9.30%, 7.60% and 7.34%, respectively. Second, CrowdSat-Net is less sensitive to the  
 653 selection of  $\gamma$  than the other methods. For example, when  $\gamma$  changes from 2.0 to 0.0 pixels, the F1-score and  
 654 Precision of FIDTMCL decrease from 64.41% to 17.59% and from 70.80% to 19.04%, respectively. P2PNet  
 655 and APGCC also show a similar trend. These steep drops illustrate that many of their predicted points lie near  
 656 the matching point and are only counted as true positives when a relatively large  $\gamma$  is used.

#### 657

#### 658 *4.5. Computational Efficiency and Memory Usage*

659

660 This section assesses the practical applicability of different CD methods by comparing the computational  
 661 cost per  $256 \times 256$  pixel patch in terms of FLOPs, the number of parameters, runtime latency and peak GPU  
 662 memory. As presented in Table 5, CrowdSat-Net is the most computationally demanding model. It has the  
 663 largest FLOPs and parameter count, and its latency and peak memory are higher than those of the other  
 664 methods. This overhead mainly stems from the Hourglass structure itself and the additional DCPAN and  
 665 HFGDU modules. Despite this larger computational cost, CrowdSat-Net yields notable localization  
 666 performance on CrowdSat, as discussed in Section 4.4. Moreover, in terms of practical deployment, for  
 667 scenarios where computational resources are very limited and slightly lower localization accuracy is  
 668 acceptable, methods such as SCALNET or APGCC may be preferred. For applications that require accurate  
 669 localization of extremely small objects under strict matching criteria, CrowdSat-Net offers a more appropriate  
 670 solution.

671

672 Table 5 Computational efficiency and memory usage for different CD methods (per  $256 \times 256$  pixel patch).

Methods	FLOPs	Params	Runtime Latency	Peak GPU memory
---------	-------	--------	-----------------	-----------------

SCALNET (Wang et al., 2021)	7.32G	18.64M	1.96ms	94.45M
P2PNet (Song et al., 2021)	26.18G	19.22M	2.21ms	122.99MB
CLTR (Liang et al., 2022b)	16.03G	59.89M	14.51ms	273.23MB
PET (Liu et al., 2023a)	221.06G	52.98M	12.24ms	300.85MB
FIDTMCL (Liang et al., 2022a)	35.60G	66.58M	15.02ms	317.24MB
APGCC (Chen et al., 2024)	23.49G	18.56M	2.039ms	113.88MB
HerdNet (Delplanque et al., 2023)	7.55G	18.16M	3.07ms	106.67MB
UNE (Wu et al., 2023)	40.19G	17.26M	2.87ms	182.42MB
CrowdSat-Net	297.29G	94.12M	20.25ms	701.05MB

673

674 *4.6. Performance of Different Backbones in CrowdSat-Net*

675

676 To investigate the performance of different backbones in CrowdSat-Net, we integrated DCPAN and  
677 HFGDU into three representative backbones, namely, U-Net (Ronneberger et al., 2015), Deep Layer  
678 Aggregation (DLA) (Yu et al., 2018) and HRNet (Sun et al., 2019), all of which include coarse and fine spatial  
679 resolution fusion strategies. As shown in Table 6, when  $\gamma=2$ , the F1-scores range for the three backbones are  
680 59.80–62.30%, which are below those of the benchmark methods, such as UNE, FIDTMCL and APGCC  
681 (64.34–65.32%). However, as  $\gamma$  decreases, the accuracy of the three backbones becomes greater than that of  
682 the benchmark methods gradually. For example, when  $\gamma=0$ , the three backbones still retain F1-scores between  
683 35.43% and 37.36%, while UNE, FIDTMCL and APGCC fall to 22.48%, 17.59% and 24.53%, respectively.  
684 With  $\gamma$  changes, the accuracies of the three backbones are more stable than those of the benchmark methods.  
685 Furthermore, compared to the Hourglass backbone, the three backbones produce clearly smaller localization  
686 accuracy at all point-matching thresholds, but they have fewer parameters than the Hourglass backbone, which  
687 makes them suitable replacements for Hourglass in applications that operate with more tolerant  
688 point-matching thresholds.

689

690

Table 6 Comparison between different backbones in CrowdSat-Net.

Methods	$\gamma$	Params	F1-score (%)	Recall (%)	Precision (%)	MLE	MAE	RMSE
U-Net (Ronneberger et al., 2015)	0	17.79M	35.43	38.09	33.13	0.00	16.56	33.77
	1		55.29	59.43	51.69	0.36		
	1.5		60.24	64.75	56.32	0.46		
	2		60.39	64.92	56.46	0.47		
DLA (Yu et al., 2018)	0	17.15M	35.51	37.63	33.61	0.00	15.59	33.54
	1		54.96	58.25	52.02	0.35		
	1.5		59.72	63.30	56.53	0.50		
	2		59.80	63.38	56.61	0.46		
HRNet (Sun et al., 2019)	0	10.16M	37.36	36.88	37.85	0.00	13.91	34.45
	1		57.50	56.77	58.26	0.35		
	1.5		60.95	60.17	61.75	0.41		
	2		62.30	61.50	63.11	0.40		

691

692 *4.7. Performance of CrowdSat-Net with Different Crowd Densities*

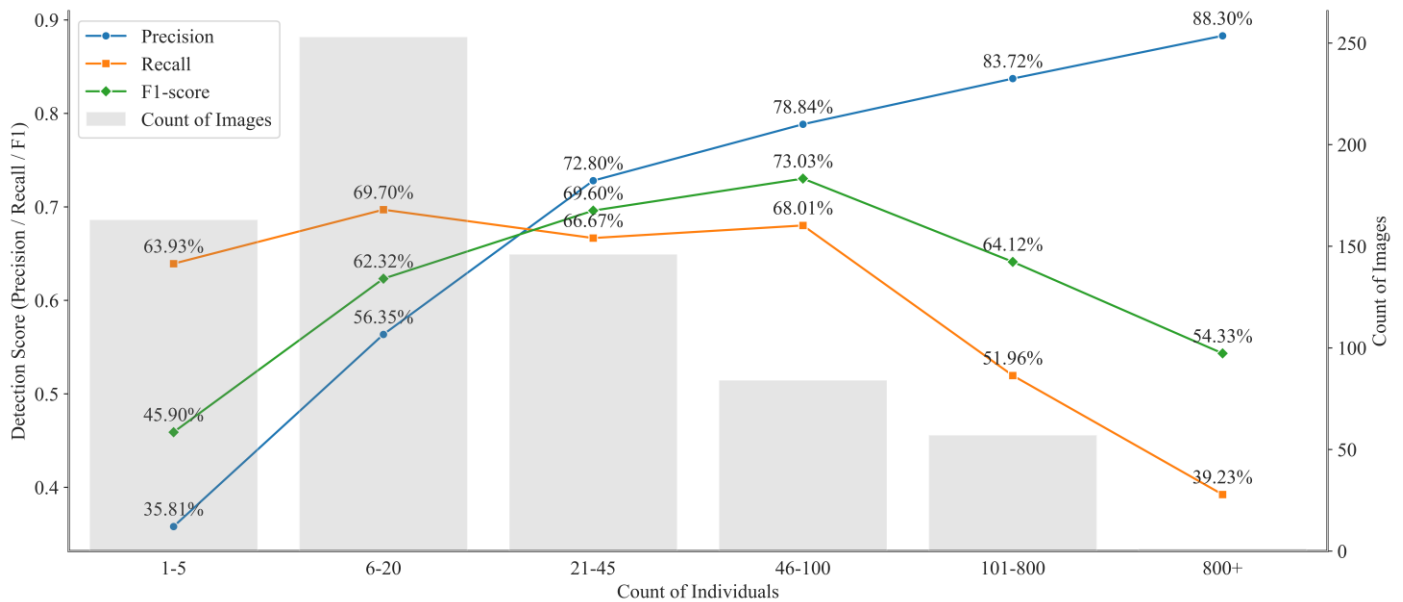
693

694 This section examines the performance of CrowdSat-Net systematically across different crowd densities to  
695 highlight its adaptive advantages in complex environments. The crowd count ranges were derived from the  
696 reference crowd count distribution using quantiles (0.1, 0.3, 0.6, 0.8, 0.95). The crowd density levels were then  
697 categorized as follows based on these quantile-derived ranges: 1–5 individuals within an image were classified  
698 as extremely sparse, 6–20 as sparse, 21–45 as moderate, 46–100 as relatively high, 101–800 as high and 800+  
699 as extremely dense.

700 As illustrated in Fig. 13, CrowdSat-Net demonstrates consistent robustness in moderate-to-high crowd  
701 density scenarios, achieving peak F1-scores of 73.03% in the 46–100 group and 69.60% in the 21–45 group.  
702 However, performance degradation was observed in both extremely sparse and extremely dense scenarios. For  
703 extremely sparse groups (1–5 individuals), Precision decreases to 35.81%, and in extremely dense crowds  
704 (800+), Recall drops to 39.23%. To explain these situations, a visual diagram of these situations is shown in

705 Fig. 14. In extremely sparse scenes, the obvious decrease in Precision is attributed to an increased  
 706 susceptibility to false positives caused by background clutter, such as image noise and other small objects (e.g.,  
 707 road asphalt, stone pillars and street lamps). Additionally, in extremely dense crowds, the decrease in Recall is  
 708 due to severe occlusion among individuals, making true positive detection challenging.

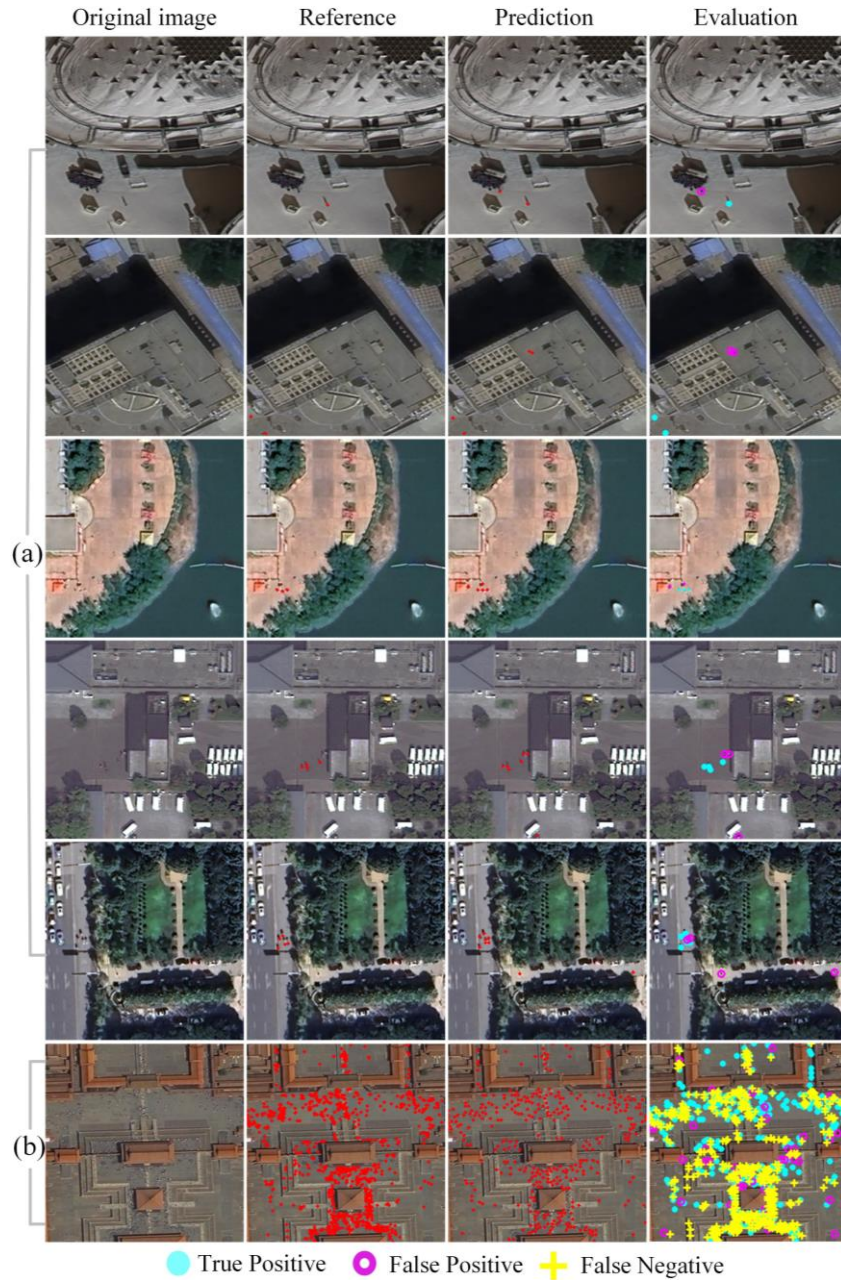
709



710

711 Fig. 13. Performance of CrowdSat-Net on CrowdSat across different crowd densities. The lines show the Precision, Recall and  
 712 F1-score evaluated on subsets of images grouped by crowd densities, while the grey bars in the background indicate the number of  
 713 images in each group.

714



715

716 Fig. 14. Examples of the localization performance of CrowdSat-Net in extreme scenarios: (a) Extremely sparse. (b) Extremely  
 717 dense.

718

#### 719 4.8. Cross-Regional Generalization of CrowdSat-Net

720

721 In this research, we used CrowdSat to cover a broad range of scenes, e.g., snowy regions, vegetated areas,  
 722 beaches, desert regions, etc., so that the generalization ability of CD models can be examined under diverse  
 723 conditions. To further assess cross-regional generalization, we additionally tested CrowdSat-Net and the

724 benchmark methods (HerdNet, UNE, SCALNet, P2PNet, FIDTMCL and APGCC) on six unseen global  
 725 regions that are not included in CrowdSat. These regions are Red Square (Moscow, Russia), Metropolitan  
 726 Cathedral (Mexico City, Mexico), Phra Nakhon (Bangkok, Thailand), Djemaa el Fna (Marrakesh, Morocco),  
 727 India Gate (New Delhi, India) and National Mall (Washington, D.C., USA). Also, these regions exhibit  
 728 distinct characteristics. Red square and Djemaa el Fna are large plazas surrounded by tall buildings, where  
 729 shadows and open areas with dark colors make individuals present low contrast. Metropolitan Cathedral and  
 730 India Gate contain complex traffic junctions, with moving vehicles mixed with pedestrians on roads. Phra  
 731 Nakhon and the National Mall include extensive vegetated areas in which dark green grass shows similar  
 732 texture characteristics to individuals. All images were annotated with the same point-based protocol as  
 733 CrowdSat. Their image sizes, collection dates, sources, spatial resolution, count of individuals and scenes  
 734 included are summarized in Table 7.

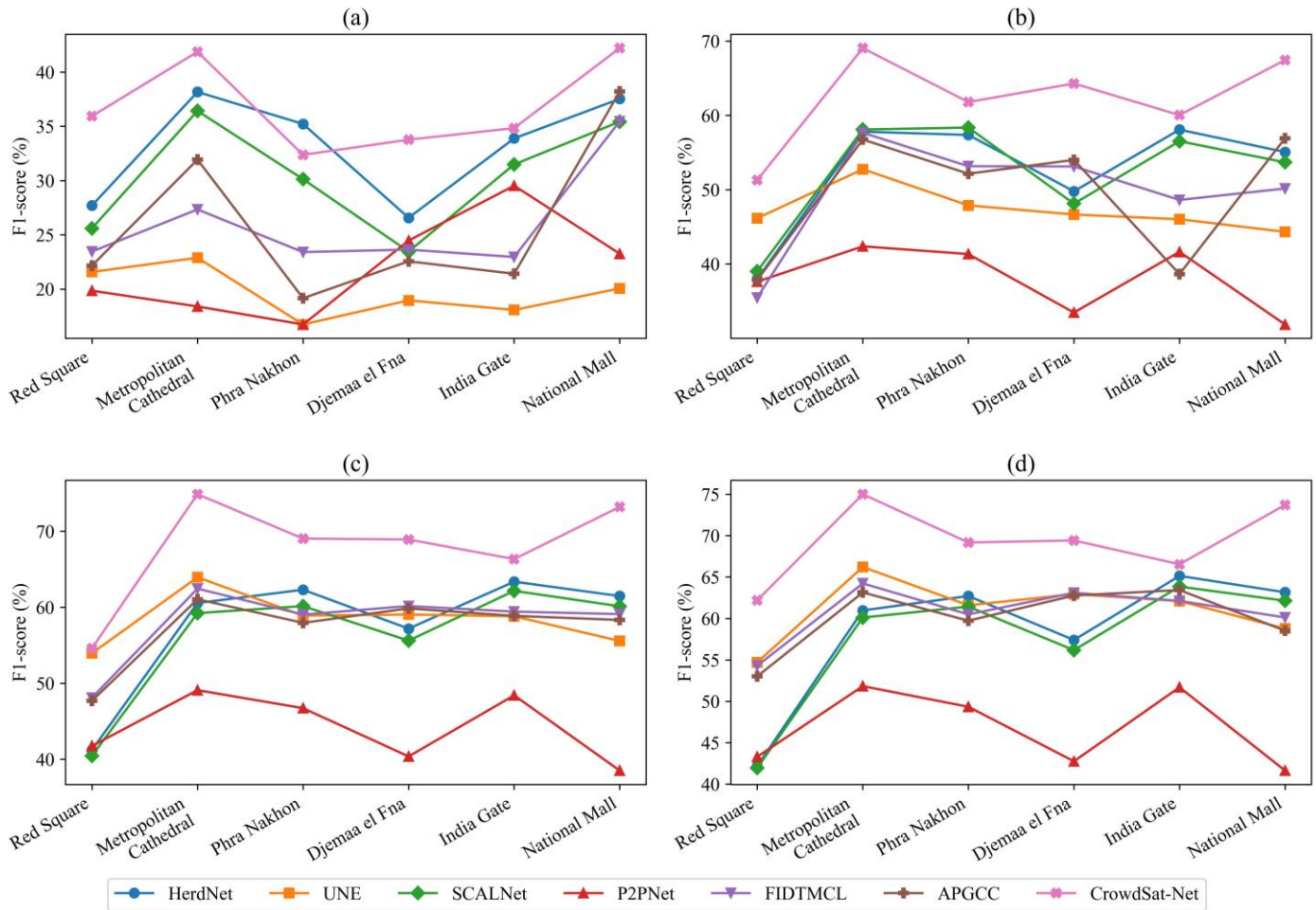
735

736

Table 7 Description of cross-regional test regions.

Regions	Image Size (pixels)	Date	Source	Spatial Resolution	Count	Scenes
Red Square	$512 \times 512$	Sep. 27, 2023	Google Earth	0.3 m	394	Open public areas; built-up areas
Metropolitan Cathedral	$1280 \times 1405$	Oct. 11, 2024	Google Earth	0.3 m	898	Built-up areas; traffic junctions
Phra Nakhon	$1402 \times 1754$	Jan. 11, 2025	Google Earth	0.3 m	725	Open public areas; vegetated areas; traffic junctions
Djemaa el Fna	$2003 \times 1514$	Mar. 14, 2024	Google Earth	0.3 m	773	Open public areas; built-up areas; traffic junctions
India Gate	$2337 \times 677$	Oct. 26, 2024	Google Earth	0.3 m	519	Open public areas; vegetated areas; traffic junctions
National Mall	$2718 \times 2079$	Mar. 3, 2024	Google Earth	0.3 m	575	Vegetated areas; traffic junctions

737



738

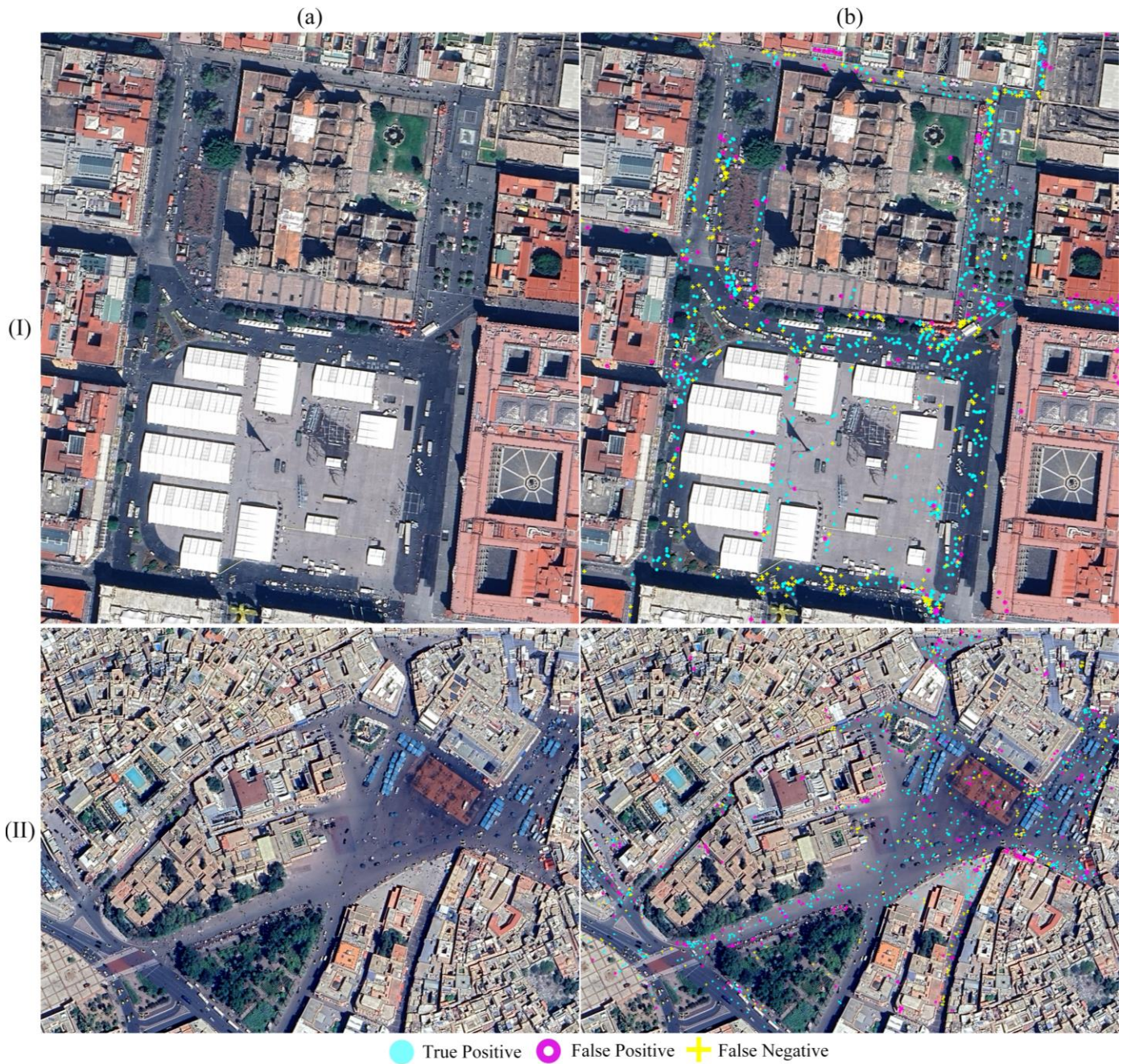
739 Fig. 15. F1-scores of different CD methods for the unseen regions under varying point-matching thresholds. (a)  $\gamma=0$ . (b)  $\gamma=1$ . (c)740  $\gamma=1.5$ . (d)  $\gamma=2$ .

741

742 As displayed in Fig. 15, across all six unseen regions, CrowdSat-Net achieves the largest or second-largest  
 743 F1-score for all point-matching thresholds. CrowdSat-Net produces mean F1-scores of 36.83%, 62.32%,  
 744 67.83% and 69.34% for  $\gamma=0.0$ , 1.0, 1.5 and 2.0, respectively. Figs. 16 and A.2 show the predictions of  
 745 CrowdSat-Net in the six unseen regions. For low-contrast plazas, CrowdSat-Net detects most individuals  
 746 successfully in open areas and along building façades. At complex traffic junctions, CrowdSat-Net is able to  
 747 separate most of the individuals from nearby vehicles. In vegetated areas, CrowdSat-Net can still localize a  
 748 large proportion of isolated individuals that are scattered across the grass. However, several obvious false  
 749 positives and negatives are also observed in these regions: false positives mostly occur on background objects  
 750 that exhibit visual characteristics similar to individuals, such as road asphalt, street lamps, stone pillars, etc..

751 False negatives are concentrated in highly dense clusters where adjacent individuals overlap. These limitations  
 752 will be further discussed in Section 5.3.

753



754

755 Fig. 16. The visual localization performance of CrowdSat-Net in unseen foreign regions. (a) Metropolitan Cathedral. (b) Djemaa el  
 756 Fna. (I) Original image. (II) Evaluation. The images in (a)-(b) were obtained from the Google Earth platform (© Google Earth 2025).

757

758

## 759 **5. Discussion**

760

### 761 *5.1. Limitations of the CrowdSat Dataset*

762

763 While CrowdSat demonstrates advantages for CD via VFR satellite imagery, several inherent limitations  
764 must be acknowledged to ensure its appropriate application and interpretation, which are as follows:

765 *1) Limited Coverage Area.* While CrowdSat was collected in regions with diverse environmental conditions,  
766 its coverage range is within China, which may limit the generalizability of models in extremely rare global  
767 environments, such as polar regions and dense tropical forests.

768 *2) Occlusion Constraints.* The reliance of CrowdSat on optical imagery limits its effectiveness inherently in  
769 occluded environments. Buildings, dense tree canopies and their shadows can obscure individuals, making  
770 crowds only partially detectable or undetectable in some scenes.

771 *3) Minimum Detectable Crowd Size.* At a spatial resolution of 0.3 m, more than one individual may fall  
772 within a single pixel in extremely high-density gatherings, but such cases are very rare in CrowdSat. To  
773 compete with current CD models, which can also not resolve sub-pixel individuals, we adopt a labelling  
774 protocol that restricts annotations to at most one individual per pixel. This strategy may lead to an  
775 underestimation of the number of individuals.

776 *4) Temporal Resolution Constraints.* Although VFR satellites offer relatively frequent revisit times (e.g.,  
777 BJ3N: every 5 days), the short-lived and time-specific nature of crowd activities makes them difficult to  
778 capture reliably. Cloud cover and atmospheric conditions often extend the effective revisit interval.

779

### 780 *5.2. Applicability of CrowdSat-Net*

781

782 In this research, CrowdSat-Net demonstrated reliable localization performance of individuals, and it is  
783 expected that CrowdSat-Net and its two innovative modules (DCPAN and HFGDO) have the potential for  
784 broader applications.

### 786 *5.2.1. Generalized Small Object Detection*

788 While CrowdSat-Net was proposed for large-scale CD, its core design supports generalized small-object  
789 detection for satellite imagery. Many non-human objects share similar attributes with human crowds: low  
790 pixel occupancy (typically 4-10 pixels per object), irregular spatial distributions and high contextual  
791 heterogeneity. Potential applications include 1) wildlife monitoring, such as detection of migratory ungulates,  
792 penguin colonies or marine species from VFR satellite imagery; 2) transportation surveillance, including  
793 identification of vehicles in road networks, parking lots or logistics hubs; and 3) environmental mapping, such  
794 as recognition of vegetation clusters or deforestation patterns in ecological studies. For example,  
795 CrowdSat-Net achieves an F1-score of 83.16%, a Precision of 83.69% and a Recall of 82.67% on the AED  
796 dataset (Naude et al., 2019), which was collected in Africa and contains 1,649 RGB images with 12,455  
797 manually annotated elephants (point annotation), showing a great localization performance in more general  
798 small-object detection tasks.

### 800 *5.2.2. Modular Integration with Existing Remote Sensing Frameworks*

802 DCPAN and HFGDO can be potentially integrated into various backbone architectures. Specifically, these  
803 two modules can be embedded into widely used coarse and fine spatial resolution fusion architectures, such as  
804 Feature Pyramid Network (FPN) (Lin et al., 2017a), U-Net (Ronneberger et al., 2015) and HRNet (Sun et al.,  
805 2019). For example, in U-Net, each standard skip connection can be embedded with DCPAN, and the  
806 traditional upsampling method (bilinear interpolation) can be replaced with HFGDO.

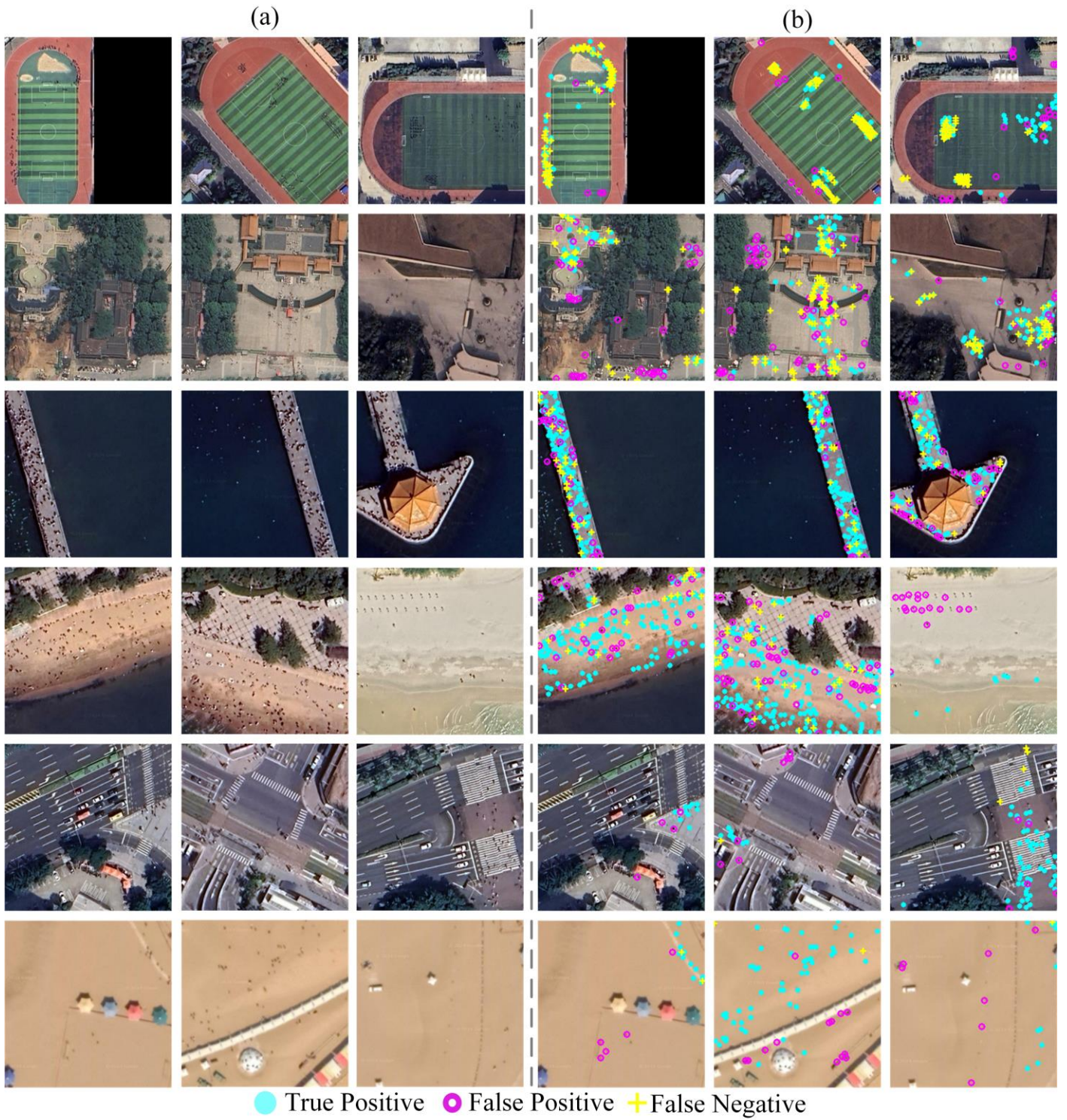
807

808 *5.3. Limitations of CrowdSat-Net*

809

810 To explore the limitations of CrowdSat-Net, we selected several representative scenes, including stadiums,  
811 parks, bridges, beaches, road junctions and desert tourist areas, where CrowdSat-Net exhibits noticeable  
812 localization errors, as shown in Fig. 17. First, a non-negligible fraction of false negatives occurs in shadows or  
813 under occlusions by trees and buildings (row 2). In these regions, CrowdSat-Net tends to miss individuals or to  
814 produce only a few sparse detections along trees or buildings. Second, in highly dense scenes, such as bridges  
815 and beaches, the signals of adjacent individuals overlap heavily, which makes CrowdSat-Net difficult to  
816 separate all of them into distinct instances and leads to false negatives within compact groups. Third, false  
817 positives are frequently observed on high-contrast background patterns, including stadium markings, beach  
818 umbrellas, E-bikes and road markings. The textures of these objects are highly similar to those of individuals,  
819 causing the network to overestimate the number of individuals in these scenes. Finally, in sparsely populated  
820 regions, isolated objects such as small kiosks or shadows of infrastructure in desert regions may also be  
821 misclassified as individuals.

822



823

824 Fig. 17. Error maps of CrowdSat-Net on representative scenes. (a) Original images. (b) Evaluation. Some of the images were  
 825 obtained from the Google Earth platform (© Google Earth 2025).

826

827 *5.4. Future Research*

828

829 Based on the limitations identified in CrowdSat and CrowdSat-Net, we propose two key directions for future  
830 research:

831 *1) Large-Scale and Multi-Source Dataset Expansion.* The CrowdSat dataset primarily covers Chinese  
832 regions, which may limit its geographical generalizability, especially for extremely rare global environments.  
833 To enhance model robustness, future datasets should incorporate global crowd patterns (by leveraging  
834 multi-source optical satellite imagery (0.3 m spatial resolution) and SAR constellations) and encompass more  
835 complex background environments (such as dense urban environments, mixed land cover regions and varied  
836 climatic conditions) to improve model adaptability across diverse scenarios.

837 *2) Temporal-Aware and Super-Resolution-Based Crowd Detection Frameworks.* To decrease the false  
838 detection rates in extremely sparse crowd regions, future models may integrate multi-temporal image  
839 sequences with dynamic change detection mechanisms. One possible solution is a dual-branch design, where  
840 one part of the model focuses on detecting individuals in an image while the other tracks how they persist over  
841 time, reducing false detections due to static, small-sized background objects. On the other hand, to alleviate the  
842 individual signal occlusion in extremely dense crowd regions, super-resolution techniques may be employed  
843 to enhance image clarity, making it easier to distinguish individuals.

844

## 845 *5.5. Ethical Considerations*

846

847 While this research focuses on the use of VFR satellite imagery for constructive and socially beneficial  
848 applications, such as urban planning, public safety monitoring, post-disaster response, etc., we acknowledge  
849 the broader discussion of the dual-use potential of CD using VFR satellite imagery. Recent research ([Ghamisi  
850 et al., 2025](#)) highlights the need to consider the privacy risks associated with fine spatial resolution remote  
851 sensing data carefully. It is important to note that this research utilizes historical satellite imagery, not  
852 real-time data. Additionally, while this research can detect the presence of individuals, the individuals are  
853 represented as indistinct black or white dot-like shapes in VFR satellite imagery, making it impossible to

854 identify them. Nevertheless, we support future research to develop privacy protection frameworks to ensure  
855 the ethical and responsible use of AI-driven geospatial analysis.

856

857

## 858 **6. Conclusion**

859

860 This research investigated CD from VFR satellite imagery and introduced both a new dataset and a tailored  
861 CD framework. Specifically, we constructed CrowdSat, a large-scale CD dataset derived from multiple VFR  
862 satellite sensors with over 120k labeled individuals across heterogeneous scenes, together with a  
863 multi-temporal background-removal strategy to reduce point-level mislabeling. Building on this dataset, we  
864 proposed CrowdSat-Net with the innovative DCPAN and HFGDU modules to enhance small-object feature  
865 representation and recover high-frequency details during upsampling. The key findings of this research are  
866 summarized as follows.

- 867 1) VFR satellite imagery overcomes the limitations of small-scale (both spatially and temporally) ground  
868 and aerial imagery, offering a promising pathway for large-scale crowd analysis in various applications.
- 869 2) CrowdSat-Net outperforms five advanced CD methods (designed for ground or aerial imagery) based  
870 on CrowdSat, achieving the largest F1-score of 66.12% and Precision of 73.23%.
- 871 3) The DCPAN and HFGDO modules are effective in increasing CD accuracy, increasing the F1-score by  
872 1.70% and Precision by 4.69%.
- 873 4) CrowdSat-Net performs reliably in moderate and relatively high crowd density scenarios, with  
874 F1-scores of 69.60% and 73.03%, respectively.
- 875 5) CrowdSat-Net demonstrates great cross-regional generalization in regions across the globe, with  
876 F1-scores of 75.00% and 73.70% in the Metropolitan Cathedral, Mexico and the National Mall, USA,  
877 respectively.

878 Overall, for large-scale CD, VFR satellite imagery offers an appropriate source, and CrowdSat-Net provides  
879 an effective solution. Future research will focus on constructing larger-scale CD datasets and developing more  
880 refined and generalized model architectures.

881

882

### 883 **CRedit Authorship Contribution Statement**

884

885 **Tong Xiao:** Conceptualization, Data curation, Formal analysis, Methodology, Validation, Visualization,  
886 Writing–original draft, Writing–review & editing. **Qunming Wang:** Conceptualization, Data curation,  
887 Funding acquisition, Supervision, Writing–original draft, Writing–review & editing. **Ping Lu:** Formal  
888 analysis, Validation, Writing–review & editing. **Tenghai Huang:** Data curation. **Xiaohua Tong:**  
889 Conceptualization, Formal analysis, Writing–review & editing. **Peter M. Atkinson:** Supervision,  
890 Conceptualization, Writing–review & editing.

891

892

### 893 **Declaration of Competing Interest**

894

895 The authors declare that they have no known competing financial interests or personal relationships that  
896 could have appeared to influence the work reported in this paper.

897

898

### 899 **Acknowledgments**

900

901 This work was supported by the National Natural Science Foundation of China under grant Nos. 42222108,  
902 42221002 and 42171345. The authors are grateful to Google Earth Platform for its satellite data support.

903

904

905 **References**

906

907 Abousamra, S., Hoai, M., Samaras, D., Chen, C., 2021. Localization in the crowd with topological constraints. In: *Proceedings of the*  
908 *AAAI Conference on Artificial Intelligence*, pp. 872–881.

909 Adam, K.D.B.J., et al., 2014. A method for stochastic optimization. *arXiv preprint* arXiv:1412.6980.

910 Al-Nami, W.T., 2023. Ranking and analysis the strategies of crowd management to reduce the risks of crushes and stampedes in  
911 crowded environments and ensure the safety of passengers. *Neutrosophic Systems with Applications* 8, 61–78.

912 Bahmanyar, R., Vig, E., Reinartz, P., 2019. MRCNet: Crowd counting and density map estimation in aerial and ground imagery.  
913 *arXiv preprint* arXiv:1909.12743.

914 Bashir S.M.A., Wang Y., 2021. Small object detection in remote sensing images with residual feature aggregation-based  
915 super-resolution and object detector network. *Remote Sensing* 13(9), 1854.

916 von Borstel, M., Kandemir, M., Schmidt, P., Rao, M.K., Rajamani, K., Hamprecht, F.A., 2016. Gaussian process density counting  
917 from weak supervision. In: *European Conference on Computer Vision*. Springer, pp. 365–380.

918 Chen, I.H., Chen, W.T., Liu, Y.W., Yang, M.H., Kuo, S.Y., 2024. Improving point-based crowd counting and localization based on  
919 auxiliary point guidance. In: *European Conference on Computer Vision*. Springer, pp. 428–444.

920 Chen, H., Shi, Z., 2020. A spatial-temporal attention-based method and a new dataset for remote sensing image change detection.  
921 *Remote Sensing* 12(10), 1662.

922 Dai, J., Qi, H., Xiong, Y., Li, Y., Zhang, G., Hu, H., Wei, Y., 2017. Deformable convolutional networks. In: *Proceedings of the IEEE*  
923 *International Conference on Computer Vision*, pp. 764–773.

924 Delplanque, A., Foucher, S., Th é au, J., Bussiere, E., Vermeulen, C., Lejeune, P., 2023. From crowd to herd counting: How to  
925 precisely detect and count African mammals using aerial imagery and deep learning? *ISPRS Journal of Photogrammetry*  
926 *and Remote Sensing* 197, 167–180.

927 Feliciani, C., Shimura, K., Nishinari, K., 2022. *Introduction to Crowd Management: Managing Crowds in the Digital Era: Theory*  
928 *and Practice*. Springer Nature.

929 Fu, J., Liu, J., Tian, H., Li, Y., Bao, Y., Fang, Z., Lu, H., 2019. Dual attention network for scene segmentation. In: *Proceedings of the*  
930 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3146–3154.

- 931 Gao, G., Liu, Q., Wang, Y., 2020a. Counting dense objects in remote sensing images. In: *ICASSP 2020-2020 IEEE International*  
932 *Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, pp. 4137–4141.
- 933 Gao, G., Liu, Q., Wang, Y., 2020b. Counting from sky: A large-scale dataset for remote sensing object counting and a benchmark  
934 method. *IEEE Transactions on Geoscience and Remote Sensing* 59, 3642–3655.
- 935 Gao, J., Han, T., Wang, Q., Yuan, Y., Li, X., 2020c. Learning independent instance maps for crowd localization. *arXiv preprint*  
936 *arXiv:2012.04164*.
- 937 Gazzawe, F., Albahar, M., 2024. Reducing traffic congestion in Makkah during Hajj through the use of AI technology. *Heliyon*  
938 10(1).
- 939 Ghamisi, P., Yu, W., Marinoni, A., Gevaert, C.M., Persello, C., Selvakumaran, S., Giroto, M., Horton, B.P., Rufin, P., Hostert, P.,  
940 Pacifici, F., Atkinson, P.M., 2025. Responsible Artificial Intelligence for Earth Observation: Achievable and realistic paths  
941 to serve the collective good. *IEEE Geoscience and Remote Sensing Magazine*, 2–26.
- 942 Guo, Y., Jiang, H., Wu, C., Zhang, L., Du, B., 2022. Density Map-based vehicle counting in remote sensing images with limited  
943 resolution. *ISPRS Journal of Photogrammetry and Remote Sensing* 189, 201–217.
- 944 Han, T., Bai, L., Liu, L., Ouyang, W., 2023. Steerer: Resolving scale variations for counting and localization via selective  
945 inheritance learning. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 21848–21859.
- 946 Hebei Provincial People’s Government, 2011. Regulations of Hebei Province on aerial surveying and photogrammetry. URL:  
947 [https://www.hebei.gov.cn/columns/67c6d89b-b5b0-4492-98d6-997dac63a29b/202309/12/0808f9e0-84c9-4580-82e9-bfe](https://www.hebei.gov.cn/columns/67c6d89b-b5b0-4492-98d6-997dac63a29b/202309/12/0808f9e0-84c9-4580-82e9-bfe3b734c875.html)  
948 [3b734c875.html](https://www.hebei.gov.cn/columns/67c6d89b-b5b0-4492-98d6-997dac63a29b/202309/12/0808f9e0-84c9-4580-82e9-bfe3b734c875.html).
- 949 Joiner, A., McFarlane, C., Rella, L., Uriarte-Ruiz, M., 2024. Problematising density: Covid-19, the crowd, and urban life. *Social &*  
950 *Cultural Geography* 25, 181–198.
- 951 Ji, S., Wei, S., Lu, M., 2018. Fully convolutional networks for multi-source building extraction from an open aerial and satellite  
952 imagery dataset. *IEEE Transactions on Geoscience and Remote Sensing* 57(1), 574–586.
- 953 Khan, M.A., Menouar, H., Hamila, R., 2023. Revisiting crowd counting: State-of-the-art, trends, and future perspectives. *Image and*  
954 *Vision Computing* 129, 104597.
- 955 Kramer, O., 2013. K-nearest neighbors. In: *Dimensionality Reduction with Unsupervised Nearest Neighbors*. Springer, pp. 13–23.
- 956 Laradji, I.H., Rostamzadeh, N., Pinheiro, P.O., Vazquez, D., Schmidt, M., 2018. Where are the blobs: Counting by localization with  
957 point supervision. In: *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 547–562.
- 958 Lempitsky, V., Zisserman, A., 2010. Learning to count objects in images. *Advances in Neural Information Processing Systems* 23.
- 959 Li, X., Zhang, J., Yang, Y., Cheng, G., Yang, K., Tong, Y., Tao, D., 2024. SFNet: Faster and accurate semantic segmentation via  
960 semantic flow. *International Journal of Computer Vision* 132, 466–489.

- 961 Lian, D., Li, J., Zheng, J., Luo, W., Gao, S., 2019. Density map regression guided detection network for RGB-D crowd counting and  
962 localization. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1821–1830.
- 963 Liang, D., Xu, W., Zhu, Y., Zhou, Y., 2022a. Focal inverse distance transform maps for crowd localization. *IEEE Transactions on*  
964 *Multimedia* 25, 6040–6052.
- 965 Liang, D., Xu, W., Bai, X., 2022b. An end-to-end transformer model for crowd localization, in: *European Conference on Computer*  
966 *Vision*. Springer, pp. 38–54.
- 967 Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S., 2017a. Feature pyramid networks for object detection. In:  
968 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2117–2125.
- 969 Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P., 2017b. Focal loss for dense object detection. In: *Proceedings of the IEEE*  
970 *International Conference on Computer Vision*, pp. 2980–2988.
- 971 Liu, C., Lu, H., Cao, Z., Liu, T., 2023a. Point-query quadtree for crowd counting, localization, and more. In: *Proceedings of the*  
972 *IEEE/CVF International Conference on Computer Vision*, pp. 1676–1685.
- 973 Liu, C., Weng, X., Mu, Y., 2019. Recurrent attentive zooming for joint crowd counting and precise localization. In: *Proceedings of*  
974 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1217–1226.
- 975 Liu, W., Lu, H., Fu, H., Cao, Z., 2023b. Learning to upsample by learning to sample. In: *Proceedings of the IEEE/CVF International*  
976 *Conference on Computer Vision*, pp. 6027–6037.
- 977 Liu, Y., Sun, P., Wergeles, N., Shang, Y., 2021. A survey and performance evaluation of deep learning methods for small object  
978 detection. *Expert Systems with Applications* 172, 114602.
- 979 Li, Z., He, W., Cheng, M., Hu, J., Yang, G., Zhang, H., 2023. SinoLC-1: the first 1 m resolution national-scale land-cover map of  
980 China created with a deep learning framework and open-access data. *Earth System Science Data* 15(11), 4749–4780.
- 981 Long, J., Shelhamer, E., Darrell, T., 2015. Fully convolutional networks for semantic segmentation. In: *Proceedings of the IEEE*  
982 *Conference on Computer Vision and Pattern Recognition*, pp. 3431–3440.
- 983 Luo, Z., Marchi, L., Chen, F., Zhang, Y., Gaspari, J., 2025. Correlating urban spatial form and crowd spatiotemporal behavior: A  
984 case study of Lhasa, China. *Cities* 160, 105812.
- 985 Meynberg, O., Cui, S., Reinartz, P., 2016. Detection of high-density crowds in aerial images using texture classification. *Remote*  
986 *Sensing* 8, 470.
- 987 Mliki, H., Arous, O., Hammami, M., 2019. Abnormal crowd density estimation in aerial images. *Journal of Electronic Imaging* 28,  
988 013047.
- 989 Mo, H., Zhang, X., Tan, J., Yang, C., Gu, Q., Hang, B., Ren, W., 2024. Countformer: Multi-view crowd counting transformer. In:  
990 *European Conference on Computer Vision*. Springer, pp. 20–40.

- 991 Naude, J., Joubert, D., 2019. The aerial elephant dataset: A new public benchmark for aerial object detection. In: *Proceedings of the*  
992 *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pp. 48-55.
- 993 Newell, A., Yang, K., Deng, J., 2016. Stacked hourglass networks for human pose estimation. In: *European Conference on*  
994 *Computer Vision*. Springer, pp. 483–499.
- 995 Ni, W., Li, H., Wang, C., Chen, Y., Lin, D., Tao, S., Xi, X., Wang, Y., Hu, J., Li, X., 2024. Forest height extraction using GF-7 very  
996 high-resolution stereoscopic imagery and Google Earth multi-temporal historical imagery. *Journal of Remote Sensing* 4,  
997 0158.
- 998 Pokhrel, S., Chhetri, R., 2021. A literature review on impact of covid-19 pandemic on teaching and learning. *Higher Education for*  
999 *the Future* 8, 133–141.
- 1000 Ren, S., He, K., Girshick, R., Sun, J., 2015. Faster R-CNN: Towards real-time object detection with region proposal networks.  
1001 *Advances in Neural Information Processing Systems* 28.
- 1002 Rezaee, K., Rezakhani, S.M., Khosravi, M.R., Moghimi, M.K., 2024. A survey on deep learning-based real-time crowd anomaly  
1003 detection for secure distributed video surveillance. *Personal and Ubiquitous Computing* 28, 135–151.
- 1004 Ronneberger, O., Fischer, P., Brox, T., 2015. U-Net: Convolutional networks for biomedical image segmentation. In: *International*  
1005 *Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, pp. 234–241.
- 1006 Rumora, L., Gašparović, M., Miler, M., Medak, D., 2020. Quality assessment of fusing Sentinel-2 and WorldView-4 imagery on  
1007 Sentinel-2 spectral band values: A case study of Zagreb, Croatia. *International Journal of Image and Data Fusion* 11, 77–  
1008 96.
- 1009 Sam, D.B., Peri, S.V., Sundararaman, M.N., Kamath, A., Babu, R.V., 2021. Locate, size, and count: Accurately resolving people in  
1010 dense crowds via detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 43(8), 2739–2751.
- 1011 Shi, W., Caballero, J., Huszár, F., Totz, J., Aitken, A.P., Bishop, R., Rueckert, D., Wang, Z., 2016. Real-time single image and video  
1012 super-resolution using an efficient sub-pixel convolutional neural network. In: *IEEE Conference on Computer Vision and*  
1013 *Pattern Recognition (CVPR)*, pp. 1874–1883.
- 1014 Shi, Q., Liu, M., Marinoni, A., Liu, X., 2023. UGS-1m: fine-grained urban green space mapping of 31 major cities in China based on  
1015 the deep learning framework. *Earth System Science Data* 15(2), 555–577.
- 1016 Shu, W., Wan, J., Tan, K.C., Kwong, S., Chan, A.B., 2022. Crowd counting in the frequency domain. In: *Proceedings of the*  
1017 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19618–19627.
- 1018 Song, Q., Wang, C., Jiang, Z., Wang, Y., Tai, Y., Wang, C., Li, J., Huang, F., Wu, Y., 2021. Rethinking counting and localization in  
1019 crowds: A purely point-based framework. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*,  
1020 pp. 3365–3374.

- 1021 Sun, K., Xiao, B., Liu, D., Wang, J., 2019. Deep high-resolution representation learning for human pose estimation. In: *Proceedings*  
1022 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5693–5703.
- 1023 Tong, K., Wu, Y., 2022. Deep learning-based detection from the perspective of small or tiny objects: A survey. *Image and Vision*  
1024 *Computing* 123, 104471.
- 1025 Wan, J., Liu, Z., Chan, A.B., 2021. A generalized loss function for crowd counting and localization. In: *Proceedings of the*  
1026 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1974–1983.
- 1027 Wang, F., Jiang, M., Qian, C., Yang, S., Li, C., Zhang, H., Wang, X., Tang, X., 2017. Residual attention network for image  
1028 classification. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3156–3164.
- 1029 Wang, J., Chen, K., Xu, R., Liu, Z., Loy, C.C., Lin, D., 2019. CARAFE: Content-aware reassembly of features. In: *Proceedings of*  
1030 *the IEEE/CVF International Conference on Computer Vision*, pp. 3007–3016.
- 1031 Wang, Q., Gao, J., Lin, W., Li, X., 2020a. NWPU-Crowd: A large-scale benchmark for crowd counting and localization. *IEEE*  
1032 *Transactions on Pattern Analysis and Machine Intelligence* 43, 2141–2149.
- 1033 Wang, Q., Shi, W., Atkinson, P.M., Zhao, Y., 2015. Downscaling MODIS images with area-to-point regression kriging. *Remote*  
1034 *Sensing of Environment* 166, 191–204.
- 1035 Wang, Q., Shi, W., Li, Z., Atkinson, P.M., 2016. Fusion of Sentinel-2 images. *Remote Sensing of Environment* 187, 241–252.
- 1036 Wang, Q., Wu, B., Zhu, P., Li, P., Zuo, W., Hu, Q., 2020b. ECA-Net: Efficient channel attention for deep convolutional neural  
1037 networks. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11534–11542.
- 1038 Wang, Y., Hou, X., Chau, L.P., 2021. Dense point prediction: A simple baseline for crowd counting and localization. In: *2021 IEEE*  
1039 *International Conference on Multimedia & Expo Workshops (ICMEW)*. IEEE, pp. 1–6.
- 1040 Wei, W., Cheng, Y., He, J., Zhu, X., 2024. A review of small object detection based on deep learning. *Neural Computing and*  
1041 *Applications* 36, 6283–6303.
- 1042 Weng, W., Wang, J., Shen, L., Song, Y., 2023. Review of analyses on crowd-gathering risk and its evaluation methods. *Journal of*  
1043 *Safety Science and Resilience* 4, 93–107.
- 1044 Wu, Z., Zhang, C., Gu, X., Duporge, I., Hughey, L. F., Stabach, J. A., Skidmore, A. K., Hopcraft, J. G. C., Lee, S. J., Atkinson, P. M.,  
1045 McCauley, D. J., Lamprey, R. H., Ngene, S. M., Wang, T., Wang, J., 2023. Deep learning enables satellite-based  
1046 monitoring of large populations of terrestrial mammals across heterogeneous landscapes. *Nature Communications* 14,  
1047 3072.
- 1048 Xu, C., Liang, D., Xu, Y., Bai, S., Zhan, W., Bai, X., Tomizuka, M., 2022. AutoScale: Learning to scale for crowd counting.  
1049 *International Journal of Computer Vision* 130, 405–434.

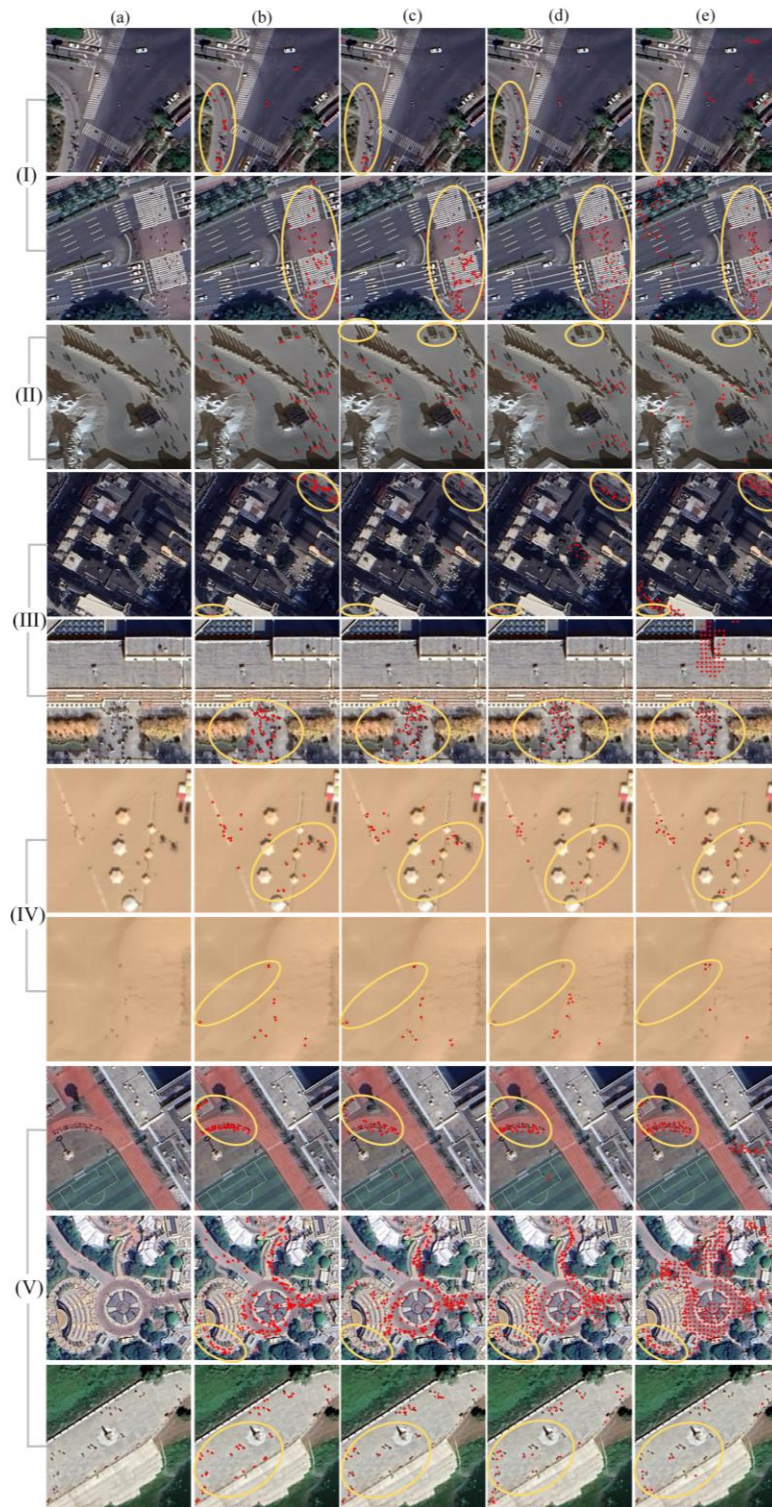
- 1050 Yi, J., Pang, Y., Zhou, W., Zhao, M., Zheng, F., 2023a. A perspective-embedded scale-selection network for crowd counting in  
1051 public transportation. *IEEE Transactions on Intelligent Transportation Systems* 25, 3420–3432.
- 1052 Yi, J., Shen, Z., Chen, F., Zhao, Y., Xiao, S., Zhou, W., 2023b. A lightweight multiscale feature fusion network for remote sensing  
1053 object counting. *IEEE Transactions on Geoscience and Remote Sensing* 61, 1–13.
- 1054 Yuan, X., Chakravarty, A., Gu, L., Wei, Z., Lichtenberg, E., Chen, T., 2025. An empirical study of methods for small object  
1055 detection from satellite imagery. *arXiv preprint arXiv:2502.03674*.
- 1056 Yun, S., Han, D., Oh, S.J., Chun, S., Choe, J., Yoo, Y., 2019. CutMix: Regularization strategy to train strong classifiers with  
1057 localizable features. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6023–6032.
- 1058 Zhang, Q., Lin, W., Chan, A.B., 2021. Cross-view cross-scene multi-view crowd counting. In: *Proceedings of the IEEE/CVF*  
1059 *Conference on Computer Vision and Pattern Recognition*, pp. 557–567.
- 1060 Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J., 2017. Pyramid Scene Parsing Network. In: *IEEE Conference on Computer Vision and*  
1061 *Pattern Recognition (CVPR)*, pp. 6230–6239.
- 1062 Zhou, X., Koltun, V., Krähenbühl, P., 2020. Tracking objects as points. In: *European Conference on Computer Vision*. Springer, pp.  
1063 474–490.
- 1064 Zhou, X., Wang, D., Krähenbühl, P., 2019. Objects as points. *arXiv preprint arXiv:1904.07850*.

1065

1066

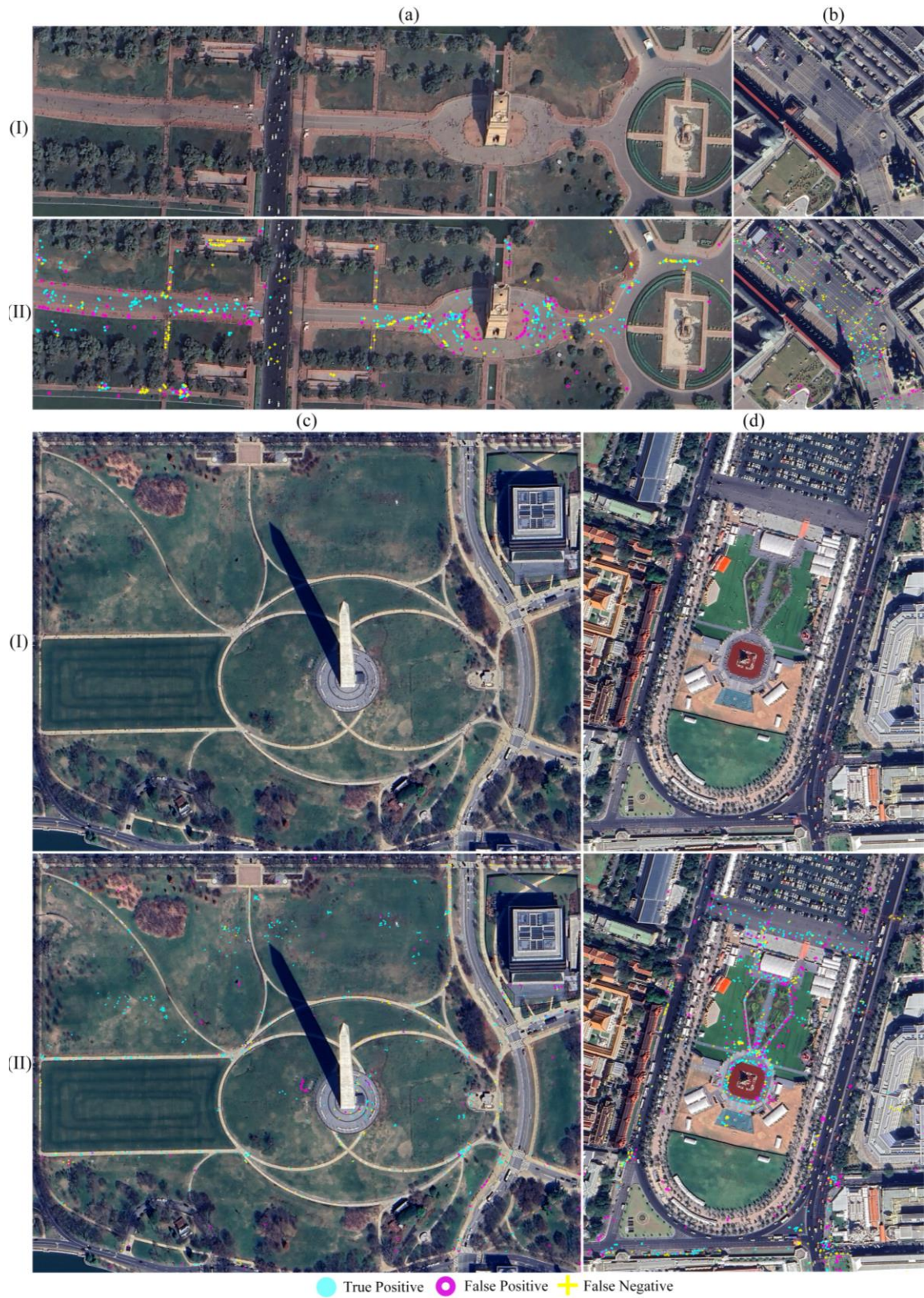
## 1067 **Appendix**

1068



1069

1070 Fig. A.1 Visual comparisons between representative CD methods (P2PNet, CLTR and PET) for the CrowdSat dataset. (a) Original  
 1071 image. (b) Reference. (c) P2PNet. (d) CLTR. (e) PET. (I) Traffic junctions. (II) Snowfields. (III) Dense urban regions. (IV) Desert  
 1072 regions. (V) Other common impervious regions. Some of the images in panels (I)-(V) were obtained from the Google Earth platform  
 1073 (© Google Earth 2025).



1074

1075 Fig. A.2. The visual localization performance of CrowdSat-Net in other unseen foreign regions. (a) India Gate. (b) Red Square. (c)  
 1076 National Mall. (d) Phra Nakhon. (I) Original image. (II) Evaluation. The images in (a)-(d) were obtained from the Google Earth  
 1077 platform (© Google Earth 2025).