

# SAFE-Net: Scale-Aware Feature Enhancement for Aerial Person Detection in Flood Disaster Imagery

Anonymous CVPR submission

Paper ID 5

## Abstract

001 *Detecting stranded persons from Unmanned Aerial Vehicles*  
002 *(UAVs) is crucial for autonomous search and rescue operations*  
003 *during flood disasters. However, reliable detection in aerial flood*  
004 *imagery remains challenging due to extreme scale variations,*  
005 *strong water reflections and motion blur. In many cases,*  
006 *human targets occupy less than 0.1% of the image area,*  
007 *making them difficult to detect using conventional object*  
008 *detectors. To address these challenges, we propose SAFE-Net*  
009 *(Scale-Aware Feature Enhancement Network), a lightweight*  
010 *detection framework which builds upon the YOLOv8*  
011 *architecture by replacing standard C2f blocks with a*  
012 *Scale-Aware Feature Enhancement (SAFE) module. The SAFE*  
013 *module improves the representation of tiny human*  
014 *targets through two mechanisms: scale-aware spatial*  
015 *weighting to emphasize extremely small objects, and*  
016 *texture enhancement using depthwise convolutions to*  
017 *recover fine edge information while suppressing noise from*  
018 *water surfaces. We also introduce UAV-based Flood Survivor*  
019 *Detection Dataset (UAV-SURV), a dataset of 6,122*  
020 *annotated aerial images collected from real flood monitoring*  
021 *videos. Experimental results show that SAFE-Net improves*  
022 *detection accuracy by 10.9% in mAP@0.5:0.95 over*  
023 *YOLOv8n while reducing model parameters by 33.8%,*  
024 *demonstrating an effective and lightweight solution for*  
025 *UAV-assisted disaster response.*

## 026 1. Introduction

027 Flood disasters are among the most destructive natural hazards  
028 worldwide, often leaving individuals stranded in inaccessible  
029 urban and rural regions. Rapid identification of survivors is  
030 critical during search-and-rescue operations, where early  
031 detection significantly improves survival rates. In recent  
032 years, Unmanned Aerial Vehicles (UAVs) have emerged as  
033 effective tools for disaster response [1, 3, 14], providing  
034 rapid aerial reconnaissance and real-time situational  
035 awareness over large flooded areas. However, man-



Figure 1. Examples of aerial flood scenes captured by UAVs during floods. The images illustrate the complexity of flooded environments, including water reflections, partially submerged structures, debris, and extremely small human instances.

ually inspecting large volumes of drone footage is time-consuming, cognitively demanding, and prone to human error. Automated computer vision systems capable of detecting stranded individuals from aerial imagery can therefore play a vital role in accelerating rescue operations.

Detecting stranded persons in aerial flood imagery remains challenging due to the complex visual characteristics of flooded environments. As shown in Figure 1, flood scenes often contain strong reflections, ripples, glare, shadows, and floating debris that introduce severe visual noise. In addition, due to the altitude of UAV platforms, human targets frequently appear extremely small, sometimes occupying less than 0.1% of the total image area. These conditions significantly degrade the performance of conventional object detectors, which are typically trained on ground-level datasets containing larger and clearer objects [2, 4, 6, 9].

Recent deep learning-based object detectors [4, 12], particularly single-stage architectures such as YOLO, have demonstrated strong performance in real-time detection tasks. However, these models still struggle in flood scenarios due to extreme scale variations, low contrast, and

057 complex background clutter. Moreover, publicly available  
058 datasets rarely capture the unique visual conditions present  
059 during real flood disasters, limiting the ability of detec-  
060 tion models to generalize effectively in such environments.  
061 These challenges highlight the need for specialized datasets  
062 and detection architectures capable of handling extremely  
063 small targets and suppressing water-induced visual noise in  
064 aerial flood imagery.

065 To address these challenges, we introduce *UAV-SURV*,  
066 a dataset consisting of 6,122 annotated aerial images ex-  
067 tracted from real flood monitoring videos. The dataset cap-  
068 tures diverse flood conditions, including variations in illu-  
069 mination, water patterns, camera altitude, and object scale.  
070 Each image is annotated with bounding boxes for human  
071 instances and divided into training and testing subsets. Fur-  
072 thermore, we propose SAFE-Net (Scale-Aware Feature En-  
073 hancement Network), a modified YOLOv8 architecture de-  
074 signed to improve aerial person detection in flood disaster  
075 scenes. SAFE-Net replaces the standard C2f modules in the  
076 YOLOv8 backbone with a Scale-Aware Feature Enhance-  
077 ment (SAFE) module. The SAFE module enhances fea-  
078 ture representation through scale-aware spatial weighting  
079 and texture enhancement mechanisms that emphasize weak  
080 human features while suppressing noise caused by water re-  
081 flections and clutter.

082 The main contributions of this work are summarized as  
083 follows:

- 084 • We introduce *UAV-SURV*, a dataset of 6,122 annotated  
085 aerial images extracted from real flood monitoring videos,  
086 designed to support research on human detection in flood  
087 disaster environments.
- 088 • We propose SAFE-Net, a modified YOLOv8 architec-  
089 ture that incorporates a Scale-Aware Feature Enhance-  
090 ment module to improve feature representation for ex-  
091 tremely small human instances.
- 092 • The proposed SAFE module enhances detection perfor-  
093 mance through scale-aware spatial weighting and texture  
094 enhancement mechanisms that suppress water-induced  
095 visual noise.
- 096 • Experimental results show that SAFE-Net improves de-  
097 tection accuracy by 10.9% in mAP@0.5:0.95 compared  
098 to YOLOv8n while reducing model parameters from  
099 3.2M to 2.12M.

## 100 2. Related Work

101 Unmanned Aerial Vehicles (UAVs) have become increas-  
102 ingly important for disaster monitoring and emergency re-  
103 sponse due to their ability to capture high-resolution im-  
104 agery from difficult-to-access regions. UAV-based vision  
105 systems enable rapid assessment of affected areas and assist  
106 in locating stranded individuals during search-and-rescue  
107 operations. Several works have explored deep learning ap-  
108 proaches for analyzing UAV imagery in disaster scenar-

ios. Ijaz et al. [5] proposed a UAV-assisted edge computing  
109 framework that deploys compressed convolutional neural  
110 networks on embedded devices for real-time disaster clas-  
111 sification. Similarly, Bashir et al. [1] introduced an effi-  
112 cient CNN-based model for disaster event classification us-  
113 ing UAV imagery, demonstrating that deep learning models  
114 can effectively analyze aerial scenes captured during emer-  
115 gencies. Although these works demonstrate the usefulness  
116 of UAV imagery for disaster analysis, they primarily focus  
117 on scene-level classification rather than detecting individual  
118 victims in complex environments. 119

120 The performance of deep learning models for disaster  
121 analysis largely depends on the availability of high-  
122 quality annotated datasets. Several datasets have been in-  
123 troduced to support flood monitoring and disaster scene  
124 understanding. FloodNet [8] provides high-resolution UAV  
125 imagery collected after Hurricane Harvey and supports mul-  
126 tiple computer vision tasks such as image classification, se-  
127 mantic segmentation, and visual question answering. While  
128 FloodNet enables comprehensive flood scene understand-  
129 ing, it does not specifically address the detection of peo-  
130 ple in flood environments. More recently, the DeepFlood  
131 dataset [3] was introduced to support accurate flood map-  
132 ping and segmentation using high-resolution aerial imagery  
133 with detailed annotations of inundated vegetation and flood-  
134 affected areas. Similarly, UrbanSARFloods [14] provides  
135 a benchmark dataset for large-scale flood mapping using  
136 Sentinel-1 SAR imagery, enabling research on flood de-  
137 tection across urban and open-area environments. Despite  
138 these advances, existing flood datasets mainly focus on  
139 flood extent estimation and segmentation tasks rather than  
140 identifying stranded individuals in aerial flood imagery. Be-  
141 yond disaster-specific datasets, several large-scale UAV per-  
142 ception datasets have been proposed to support aerial scene  
143 understanding. For example, UAVScenes [13] introduces a  
144 multi-modal UAV dataset that includes synchronized cam-  
145 era images and LiDAR data with semantic annotations to  
146 enable tasks such as segmentation, depth estimation, and lo-  
147 calization. While such datasets advance UAV perception re-  
148 search, they are primarily designed for general aerial scene  
149 understanding and do not focus on disaster response or vic-  
150 tim detection scenarios.

151 Remote sensing imagery presents several challenges  
152 compared to conventional computer vision datasets. Ob-  
153 jects often appear at very small scales due to high imaging  
154 altitude, and scenes contain complex background structures  
155 with clutter, occlusions, and illumination variations. In ad-  
156 dition, obtaining pixel-level annotations for remote sens-  
157 ing images is time-consuming and expensive. To address  
158 limited labeled data, Lu et al. [7] proposed an uncertainty-  
159 aware semi-supervised learning framework that improves  
160 segmentation performance by leveraging both labeled and  
161 unlabeled remote sensing data. Despite these advances, de-

162 tecting extremely small objects such as humans in aerial  
163 flood imagery remains challenging due to scale variations,  
164 water reflections, and background clutter.

165 In contrast to existing works, this paper focuses specif-  
166 ically on detecting stranded individuals in flood scenarios  
167 using aerial imagery. We introduce the UAV-SURV dataset  
168 together with SAFE-Net, a scale-aware feature enhance-  
169 ment network designed to improve the detection of ex-  
170 tremely small human instances in complex flood environ-  
171 ments.

### 172 3. UAV-SURV Dataset Curation and Statistics

173 Reliable detection of stranded persons in aerial flood im-  
174 agery requires datasets that capture realistic flood condi-  
175 tions. However, publicly available datasets rarely contain  
176 aerial imagery recorded during real flood disasters. To ad-  
177 dress this limitation, we construct a dataset called *UAV-*  
178 *SURV*, specifically designed for aerial person detection in  
179 flood environments.

#### 180 3.1. Dataset Curation

181 Constructing a dataset for aerial person detection in flood  
182 environments is challenging due to the extremely small size  
183 of human instances, dynamic water surfaces, and the high  
184 cost of manually annotating long UAV video sequences.  
185 Similar challenges have been observed in aerial small-  
186 object detection datasets such as RealDroneVision [11],  
187 which highlights the importance of specialized datasets and  
188 semi-automatic labeling strategies. To efficiently gener-  
189 ate reliable annotations, we adopt a semi-automatic dataset  
190 creation pipeline inspired by the Self-Annotated Labeling  
191 from Videos (SA-LfV) framework [10, 11], combined with  
192 a human-in-the-loop verification process.

193 The UAV-SURV dataset is derived from drone videos  
194 recorded during real flood events in Amaravathi, India. As  
195 shown in Figure 1 these videos capture large flooded regions  
196 with complex visual conditions including water reflections,  
197 ripples, floating debris, and varying illumination. Such en-  
198 vironmental factors significantly increase the difficulty of  
199 detecting stranded individuals from aerial viewpoints.

200 The dataset construction process follows four stages:

- 201 1. Initial tracking-based proposals: Raw UAV videos were  
202 first processed using an object tracker to generate candi-  
203 date bounding boxes across frames. These tracker out-  
204 puts provided an initial set of pseudo-labels represent-  
205 ing potential human instances. Although efficient, the  
206 automatically generated proposals often contained false  
207 positives caused by background clutter or reflections and  
208 occasionally missed true human instances.
- 209 2. Human verification of tracker outputs: Human annota-  
210 tors reviewed the tracker-generated proposals to remove  
211 false positives and identify missed human instances.  
212 This step produced a cleaner initial annotation set while

Table 1. Statistical analysis of object area ratio and object aspect ratio in the dataset.

Split	Object Area Ratio			Object Aspect Ratio		
	Min	Mean	Max	Min	Mean	Max
Train	0.000042	0.000967	0.091016	0.023	0.479	14.063
Test	0.000039	0.000997	0.079982	0.062	0.478	7.031

213 significantly reducing the manual effort compared to ex-  
214 haustive frame-by-frame labeling.

- 215 3. False negative recovery via detector training: To fur-  
216 ther improve annotation coverage, a detection model was  
217 trained using the refined tracker annotations. The trained  
218 detector was then applied to the original video sequences  
219 to identify additional candidate detections that were pre-  
220 viously missed.
- 221 4. Iterative human-in-the-loop refinement: The new de-  
222 tectations generated by the trained model were again in-  
223 spected by human annotators. Remaining false positives  
224 were removed and additional missed instances were cor-  
225 rected. This iterative refinement process gradually im-  
226 proved the annotation quality while minimizing manual  
227 labeling effort.

228 Using this semi-automatic pipeline, we generated a  
229 high-quality dataset with reliable bounding box annota-  
230 tions while significantly reducing annotation cost. The fi-  
231 nal UAV-SURV dataset contains 6,122 annotated aerial im-  
232 ages representing diverse flood scenarios suitable for train-  
233 ing and evaluating aerial person detection models.

#### 234 3.2. Dataset Statistics

235 The final UAV-SURV dataset contains a total of 6,122 anno-  
236 tated images. The dataset is divided into training and testing  
237 subsets to enable reliable model evaluation, with 4,898 im-  
238 ages used for training and 1,224 images used for testing.

239 A key characteristic of aerial flood imagery is the ex-  
240 tremely small size of human instances. Due to the altitude  
241 of the UAV, individuals often occupy only a very small frac-  
242 tion of the image area. Analysis of the bounding box statis-  
243 tics shows that most human instances occupy less than 0.1%  
244 of the image area, making the detection task particularly  
245 challenging for conventional object detectors. Representa-  
246 tive samples from the UAV-SURV dataset are shown in Fig-  
247 ure 4.

248 To further analyze the dataset characteristics, we exam-  
249 ine several statistical properties of the annotations, includ-  
250 ing aspect ratio distribution, bounding box area distribution,  
251 the number of boxes per image, and the spatial distribution  
252 of human instances. Figures 2 and 3 illustrate these statis-  
253 tics for the training and testing subsets, and detailed statis-  
254 tics are provided in Table 1.

255 The aspect ratio distribution shows that most human  
256 bounding boxes fall within a narrow range, reflecting the

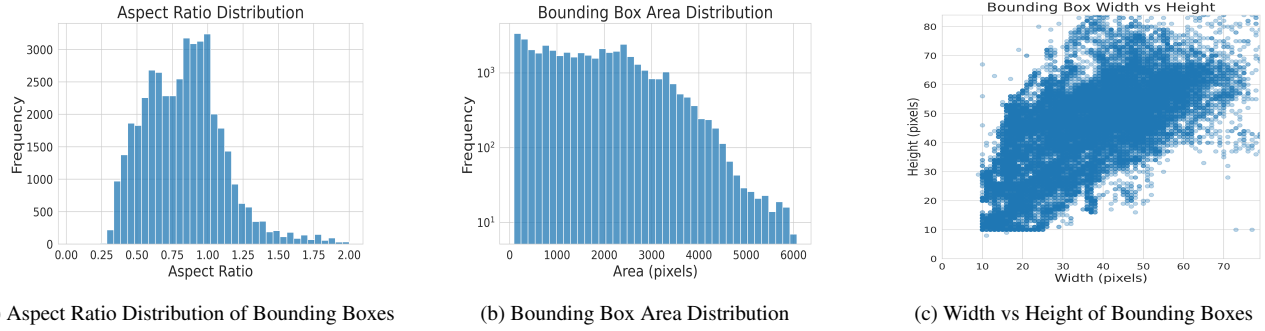


Figure 2. Statistical analysis of bounding box characteristics in the training dataset.

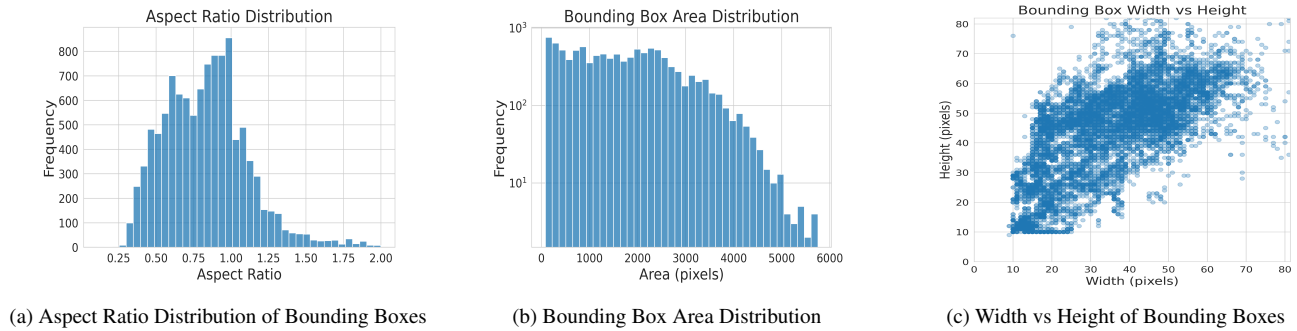


Figure 3. Statistical analysis of bounding box characteristics in the testing dataset.



Figure 4. Samples from the UAV-SURV dataset. The images illustrate aerial flood environments where stranded individuals appear at extremely small scales and diverse spatial locations. The scenes exhibit challenging visual conditions including cluttered backgrounds, water reflections, varying illumination, and complex urban layouts.

257 typical upright structure of human figures observed from  
 258 aerial viewpoints. The bounding box area distribution indicates that the majority of human instances correspond to  
 259 extremely small objects in the image. The distribution of  
 260 boxes per image shows variability in the number of persons  
 261

present in different frames, ranging from sparse scenes with  
 262 only one or two individuals to frames containing multiple  
 263 persons. These statistics confirm that UAV-SURV presents a  
 264 challenging dataset characterized by small objects, spatial  
 265 sparsity, and complex environmental conditions.  
 266

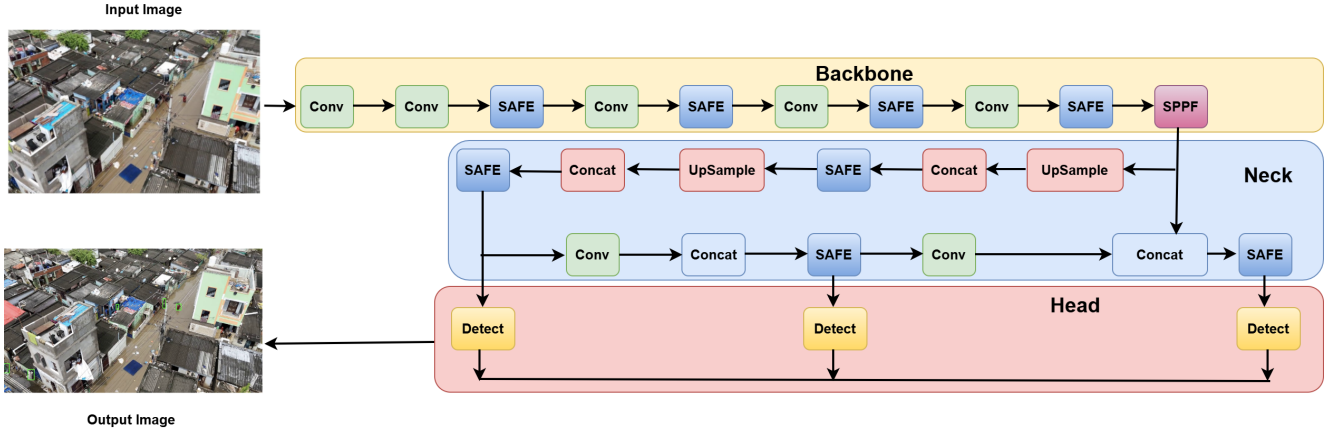


Figure 5. Overview of the proposed SAFE-Net architecture. SAFE-Net follows the YOLOv8 pipeline with backbone, neck, and detection head, where the standard C2f modules are replaced with SAFE blocks to enhance feature representations of small human instances in aerial flood imagery.

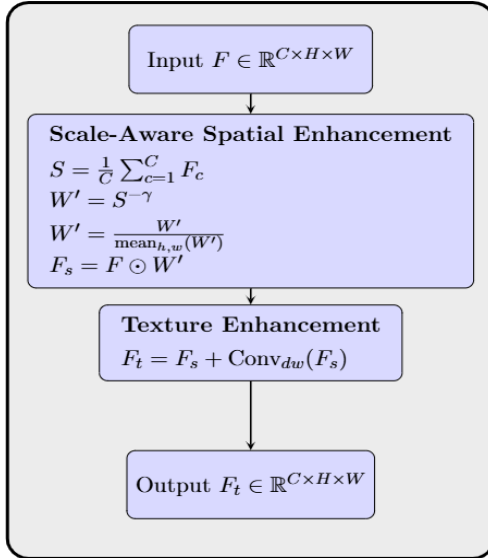


Figure 6. Overview of the proposed SAFE module used in SAFE-Net. The module refines backbone feature maps through scale-aware spatial enhancement and texture enhancement to improve the detection of small human instances in aerial flood imagery.

267 Overall, UAV-SURV provides a realistic benchmark for  
 268 aerial person detection in flood disaster environments and  
 269 supports the evaluation of detection models designed for  
 270 UAV-assisted search-and-rescue applications. The overall  
 271 SAFE-Net architecture is illustrated in Figure 5. The overall  
 272 SAFE-Net architecture is illustrated in Figure 5.

## 273 4. Proposed Methodology

274 This section presents SAFE-Net, a modified YOLOv8 ar-  
 275 chitecture designed for aerial person detection in flood dis-

276 aster environments. Detecting humans in aerial flood im-  
 277 agery is challenging due to extremely small object sizes,  
 278 water reflections, and cluttered backgrounds. To address  
 279 these challenges, SAFE-Net enhances feature representa-  
 280 tions using scale-aware spatial enhancement and texture en-  
 281 hancement mechanisms integrated into the backbone net-  
 282 work. Specifically, the standard C2f modules in YOLOv8  
 283 are replaced with a lightweight SAFE module that strength-  
 284 ens the representation of small human instances while main-  
 285 taining computational efficiency.

### 4.1. SAFE-Net Architecture 286

287 SAFE-Net follows the general YOLOv8 architecture con-  
 288 sisting of a backbone for feature extraction, a neck for  
 289 multi-scale feature aggregation, and a detection head for  
 290 bounding box prediction. The backbone extracts hierarchi-  
 291 cal feature representations from the input image, while the  
 292 neck aggregates features across multiple scales to improve  
 293 small object detection.

294 In the proposed framework, the conventional C2f mod-  
 295 ules in the YOLOv8 backbone are replaced with SAFE  
 296 blocks that integrate scale-aware spatial enhancement and  
 297 texture enhancement operations. The overall structure of  
 298 the SAFE module is illustrated in Figure 6.

299 Given an input image  $I$ , the backbone extracts feature  
 300 maps

$$F \in \mathbb{R}^{C \times H \times W} \quad (1) \quad 301$$

302 where  $C$ ,  $H$ , and  $W$  denote the number of channels,  
 303 height, and width of the feature map. These feature maps  
 304 are subsequently refined by the SAFE modules before be-  
 305 ing forwarded to the detection head.

## 306 4.2. Scale-Aware Spatial Enhancement

307 In aerial flood imagery, humans often occupy only a small  
308 fraction of the image area, leading to weak feature re-  
309 sponses that may be suppressed during feature aggregation.  
310 To address this issue, SAFE-Net introduces a scale-aware  
311 spatial enhancement mechanism that emphasizes spatial re-  
312 gions corresponding to small objects.

313 Given the feature map  $F$ , a spatial importance map is  
314 computed by averaging the feature responses across chan-  
315 nels:

$$316 \quad S = \frac{1}{C} \sum_{c=1}^C F_c \quad (2)$$

317 where  $F_c \in \mathbb{R}^{H \times W}$  denotes the  $c$ -th channel of the fea-  
318 ture map. The resulting spatial map  $S \in \mathbb{R}^{1 \times H \times W}$  captures  
319 the average activation across channels.

320 A scale-aware weighting map is then computed as

$$321 \quad W' = S^{-\gamma} \quad (3)$$

322 where  $\gamma$  controls the strength of the enhancement. To  
323 stabilize training, the weight map is normalized as

$$324 \quad W' = \frac{W'}{\text{mean}_{h,w}(W')}. \quad (4)$$

325 The normalized weight map is broadcast along the chan-  
326 nel dimension and applied to the feature map through  
327 element-wise multiplication:

$$328 \quad F_s = F \odot W' \quad (5)$$

329 where  $\odot$  denotes element-wise multiplication. This op-  
330 eration amplifies weak spatial responses associated with  
331 small-object regions while suppressing dominant back-  
332 ground responses.

## 333 4.3. Texture Enhancement

334 Flood scenes often contain low-contrast human silhou-  
335 ettes that blend with surrounding water surfaces. To  
336 enhance local structural information, SAFE-Net incorpo-  
337 rates a lightweight texture enhancement operation based on  
338 depthwise convolution.

339 The enhanced feature representation is computed as

$$340 \quad F_t = F_s + \text{Conv}_{dw}(F_s) \quad (6)$$

341 where  $\text{Conv}_{dw}$  denotes a depthwise convolution applied  
342 independently to each channel. This residual operation en-  
343 hances edge and structural information, enabling the net-  
344 work to better distinguish human shapes from water pat-  
345 terns and background clutter.

346 The resulting feature map  $F_t$  is forwarded to the  
347 YOLOv8 neck and detection head for final object predic-  
348 tion.

## 4.4. Integration with YOLOv8

The SAFE module is integrated into the YOLOv8 backbone  
by replacing the standard C2f blocks with the proposed fea-  
ture refinement block. This modification allows the detec-  
tor to better capture small human features in aerial flood  
imagery without significantly increasing model complexity.  
The refined feature maps are then passed through the neck  
and detection head for multi-scale object prediction.

By combining scale-aware spatial enhancement and tex-  
ture refinement, SAFE-Net improves the detection of small  
human instances in challenging flood environments while  
maintaining the efficiency required for UAV-based search-  
and-rescue applications.

## 5. Implementation Details

The proposed SAFE-Net model is implemented using the  
Ultralytics YOLOv8 framework in PyTorch. The stan-  
dard C2f modules in the YOLOv8 backbone are replaced  
with the proposed SAFE module, which incorporates scale-  
aware spatial enhancement and texture enhancement to im-  
prove feature representation for small human instances in  
aerial flood imagery.

### 5.1. Training Setup

The model is trained on the UAV-SURV dataset, which con-  
sists of 6,122 annotated aerial images collected from real  
flood monitoring videos. The dataset is divided into 4,898  
images for training and 1,224 images for testing. Each im-  
age contains bounding box annotations corresponding to the  
person class.

All experiments are conducted using an NVIDIA RTX  
2080 Ti GPU. The input image resolution is set to  $640 \times 640$ ,  
and the model is trained for 100 epochs with a batch  
size of 16. The training procedure follows the standard  
YOLOv8 optimization pipeline with stochastic gradient de-  
scent (SGD) optimizer, momentum of 0.937, and weight de-  
cay of  $5 \times 10^{-4}$ . Data augmentation techniques including  
mosaic augmentation, random scaling, flipping, and color  
jitter are applied during training to improve model general-  
ization.

### 5.2. Evaluation Metrics

Model performance is evaluated using standard object  
detection metrics including precision, recall, mean Av-  
erage Precision at IoU threshold 0.5 (mAP@0.5), and  
mean Average Precision across multiple IoU thresholds  
(mAP@0.5:0.95). These metrics provide a comprehensive  
assessment of detection accuracy and localization perfor-  
mance.

### 5.3. Implementation Efficiency

SAFE-Net maintains a lightweight architecture suitable for  
aerial deployment scenarios. The final model contains ap-

Table 2. Comparison of lightweight detectors on the UAV-SURV dataset. SAFE-Net achieves the best detection performance while maintaining a lightweight architecture.

Model	Params (M)	FLOPs (G)	Precision	Recall	mAP@0.5	mAP@0.5:0.95
YOLOv5n	2.6	7.7	0.963	0.807	0.880	0.710
YOLOv8n	3.2	8.7	0.967	0.813	0.889	0.725
YOLOv10n	2.3	6.7	0.962	0.804	0.889	0.740
RT-DETR	20	60	0.953	0.79	0.82	0.723
SAFE-Net (Ours)	2.12	5.9	<b>0.980</b>	<b>0.895</b>	<b>0.950</b>	<b>0.834</b>

Table 3. Ablation study of SAFE-Net components on UAV-SURV.

Method	Scale	Texture	Params (M)	FLOPs (G)	Precision	Recall	mAP@0.5:0.95
YOLOv8n + Scale Enhancement	✓	–	2.06	5.6	0.972	0.872	0.778
YOLOv8n + Texture Enhancement	–	✓	2.08	5.7	0.973	0.892	0.821
SAFE-Net (Scale + Texture)	✓	✓	2.12	5.9	<b>0.980</b>	<b>0.895</b>	<b>0.834</b>

398 proximately 2.12 million parameters and requires about 5.9  
 399 GFLOPs. During inference, the model achieves an average  
 400 inference time of approximately 3 ms per image on an  
 401 RTX 2080 Ti GPU, enabling efficient processing of aerial  
 402 imagery for disaster response applications.

## 403 6. Experiments and Results

404 This section evaluates the performance of the proposed  
 405 SAFE-Net architecture on the UAV-SURV dataset and compares  
 406 it with several lightweight YOLO-based detectors. We also  
 407 conduct an ablation study to analyze the contribution of  
 408 each component in the proposed SAFE module.

### 409 6.1. Comparison with YOLO Baselines

410 To demonstrate the effectiveness of SAFE-Net, we compare  
 411 it with three lightweight object detection models: YOLOv5n,  
 412 YOLOv8n, and YOLOv10n. All models are trained using the  
 413 same dataset split and training configuration described in  
 414 Section 5. Performance is evaluated using precision, recall,  
 415 and mean Average Precision at IoU threshold 0.5 (mAP@0.5).  
 416

417 Table 2 summarizes the detection performance of different  
 418 models on the UAV-SURV testing set. Qualitative detection  
 419 examples are shown later in Figure 8.

420 The results show that SAFE-Net consistently outperforms  
 421 all baseline detectors. In particular, SAFE-Net achieves a  
 422 10.9-point improvement in mAP@0.5:0.95 compared to  
 423 YOLOv8n while maintaining a lightweight architecture. This  
 424 improvement demonstrates the effectiveness of the proposed  
 425 SAFE module in enhancing feature representations for extremely  
 426 small human instances in aerial flood imagery, as visualized  
 427 in Figure 7.

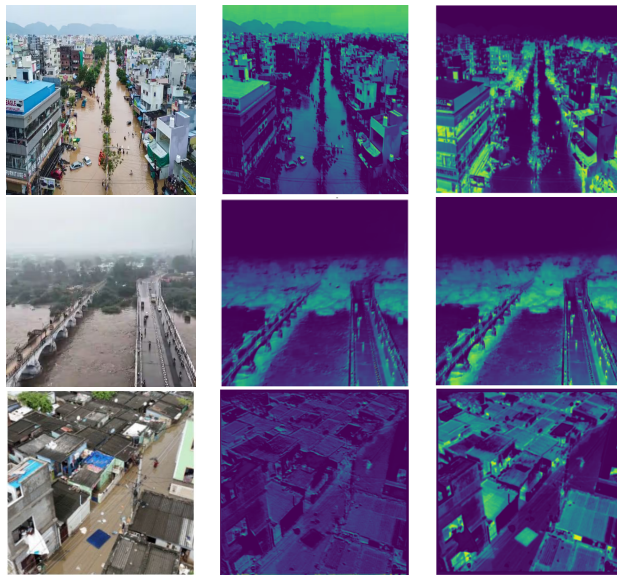


Figure 7. Visualization of feature responses produced by YOLOv8 and SAFE-Net. Each row shows the original image (left), YOLOv8n activation map (middle), and SAFE-Net activation map (right). The SAFE module produces more focused responses around human regions while suppressing background noise.

### 428 6.2. Ablation Study

429 To analyze the contribution of the implemented SAFE components,  
 430 we perform an ablation study by evaluating scale-aware enhancement  
 431 and texture enhancement. All ablation experiments are conducted  
 432 using the YOLOv8 backbone while keeping the training settings  
 433 unchanged.

434 Table 3 presents the ablation results.

435 The ablation results show that both implemented SAFE components  
 436 contribute to improved detection performance. Scale-aware  
 437 enhancement improves the detection of ex-



Figure 8. Visual results on the UAV-SURV dataset using SAFE-Net model.

438 tremely small human instances by amplifying weak spatial  
439 responses. Texture enhancement further improves boundary  
440 information, which helps distinguish human shapes from  
441 complex water textures.

442 SAFE-Net achieves the highest detection performance  
443 while maintaining a lightweight architecture. As shown  
444 in Table 2, the proposed model improves mAP@0.5 by  
445 6.1 points and mAP@0.5:0.95 by 10.9 points compared to  
446 YOLOv8n while keeping low parameter count and compu-  
447 tational complexity.

## 448 7. Conclusion

449 This work presents UAV-SURV, a real flood-disaster aerial  
450 dataset, and SAFE-Net, a lightweight YOLOv8-based de-  
451 tector with scale-aware and texture enhancement modules.  
452 Experiments show consistent gains over lightweight base-  
453 lines, particularly for tiny-person detection, while preserv-  
454 ing efficient inference suitable for UAV-assisted search-and-  
455 rescue deployment.

## 456 References

457 [1] Munzir Hubiba Bashir, Musheer Ahmad, Danish Raza Rizvi,  
458 and Ahmed A Abd El-Latif. Efficient cnn-based disaster  
459 events classification using uav-aided images for emergency  
460 response application. *Neural Computing and Applications*,  
461 36(18):10599–10612, 2024. 1, 2

462 [2] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christo-  
463 pher KI Williams, John Winn, and Andrew Zisserman. The  
464 pascal visual object classes challenge: A retrospective. *Inter-  
465 national journal of computer vision*, 111(1):98–136, 2015. 1

466 [3] Mulham Fawakherji, Jeffrey Blay, Matilda Anokye, Leila  
467 Hashemi-Beni, and Jennifer Dorton. Deepflood for inun-  
468 dated vegetation high-resolution dataset for accurate flood  
469 mapping and segmentation. *Scientific Data*, 12(1):271, 2025.  
470 1, 2

471 [4] Muhammad Hussain. Yolov1 to v8: Unveiling each variant–  
472 a comprehensive review of yolo. *IEEE access*, 12:42816–  
473 42833, 2024. 1

474 [5] Haris Ijaz, Rizwan Ahmad, Rehan Ahmed, Waqas Ahmed,  
475 Yan Kai, and Wu Jun. A uav-assisted edge framework for  
476 real-time disaster management. *IEEE Transactions on Geo-  
477 science and Remote Sensing*, 61:1–13, 2023. 2

[6] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, 478  
Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence 479  
Zitnick. Microsoft coco: Common objects in context. In 480  
*European conference on computer vision*, pages 740–755. 481  
Springer, 2014. 1 482

[7] Xiaoqiang Lu, Lingling Li, Licheng Jiao, Xu Liu, Fang Liu, 483  
Wenping Ma, and Shuyuan Yang. Uncertainty-aware semi- 484  
supervised learning segmentation for remote sensing images. 485  
*IEEE Transactions on Multimedia*, 2025. 2 486

[8] Maryam Rahnemoonfar, Tashnim Chowdhury, Argho 487  
Sarkar, Debvrat Varshney, Masoud Yari, and Robin Rober- 488  
son Murphy. Floodnet: A high resolution aerial imagery 489  
dataset for post flood scene understanding. *IEEE Access*, 490  
9:89644–89654, 2021. 2 491

[9] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 492  
Faster r-cnn: Towards real-time object detection with region 493  
proposal networks. *IEEE transactions on pattern analysis 494  
and machine intelligence*, 39(6):1137–1149, 2016. 1 495

[10] Arun Kumar Sivapuram, Prashanth Komuravelli, and Rama 496  
Krishna Sai Gorthi. Sa-1fv: self-annotated labeling from 497  
videos for object detection. *Machine Learning*, 114(1):21, 498  
2025. 3 499

[11] Arun Kumar Sivapuram, Pranav RT Peddinti, Harish Pup- 500  
pala, Komuravelli Prashanth, Jaladi Sri Harsha, and Rama 501  
Krishna Sai Gorthi. Realdronevision: Dataset and architec- 502  
ture advancements for small-object drone detection. In *Pro- 503  
ceedings of the IEEE/CVF Winter Conference on Applica- 504  
tions of Computer Vision*, pages 6687–6695, 2026. 3 505

[12] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: A 506  
simple and strong anchor-free object detector. *IEEE trans- 507  
actions on pattern analysis and machine intelligence*, 44(4): 508  
1922–1933, 2020. 1 509

[13] Sijie Wang, Siqi Li, Yawei Zhang, Shangshu Yu, Shenghai 510  
Yuan, Rui She, Quanjiang Guo, JinXuan Zheng, Ong Kang 511  
Howe, Leonrich Chandra, et al. Uavscenes: A multi-modal 512  
dataset for uavs. In *Proceedings of the IEEE/CVF Interna- 513  
tional Conference on Computer Vision*, pages 28946–28958, 514  
2025. 2 515

[14] Jie Zhao, Zhitong Xiong, and Xiao Xiang Zhu. Ur- 516  
bansarfloods: Sentinel-1 slc-based benchmark dataset for 517  
urban and open-area flood mapping. In *Proceedings of 518  
the IEEE/CVF Conference on Computer Vision and Pattern 519  
Recognition*, pages 419–429, 2024. 1, 2 520