

SIM-AIR Dataset and YOLO-KMM Model for Air-to-Air Infrared Small Target Detection

Luyi Zhang, Limin Liu, Chaowen Zheng, Haojie Yang, and Qiang Fu*

¹ Shijiazhuang Campus, Army Engineering University of PLA, Shijiazhuang 050003, PR China

Abstract – To address the lack of dedicated datasets for infrared detection of small UAVs in air-to-air scenarios, this paper first constructs the self-built SIM-AIR dataset covering complex scenarios, and then proposes YOLO-KMM an efficient YOLOv11-based object detection model tailored to the dataset's small-target characteristics and deployment requirements; collected by an UAV equipped with an infrared thermal imager, the SIM-AIR dataset consists of 3,993 precisely annotated images across four weather conditions: sunny, cloudy, snowy, and hazy, where 99.7% of the targets are ultra-small objects and their width < 40 pixels, with an average size of 11.2×6.6 pixels, including complex scenarios such as "dark targets" in snowy weather and low signal-to-noise ratio (SNR) in haze, which fully simulate real-world detection challenges. To tackle the issues of sparse small-target features and strong background interference, YOLO-KMM integrates the C2KD feature enhancement module and C3K2-MU lightweight detection head, forming a dual-optimized architecture of "feature enhancement - efficient detection": the C2KD module captures weak small-target features and suppresses noise via cross-scale fusion and attention mechanisms, while the C3K2-MU module adopts grouped convolution and depthwise separable convolution to reduce the number of parameters while preserving feature representation capability. Experiments on the SIM-AIR dataset show that YOLO-KMM achieves an mAP_{50} of 88.2%. This is 7.8 percentage points higher than the baseline YOLOv11, with a precision of 94.0% and recall of 74.3%, reduces the small-target missed detection rate by 12.5%, and maintains an inference speed of 246.18 FPS, 2.3M parameters, and 5.4 GFLOPs of computation; compared with YOLOv5/8/12, the model achieves a better balance among accuracy, speed, and complexity, verifying the practicality and challenge of the SIM-AIR dataset and providing an efficient solution for air-to-air small-target infrared detection.

Keywords. Air-to-air infrared detection; Small UAV target; SIM-AIR dataset; YOLO-KMM model; Feature enhancement; Lightweight object detection

1. Introduction

The detection of small air-to-air drones has become a crucial technological necessity in the fields of civilian security, transportation, and emergency rescue. In civil aviation, the International Civil Aviation Organization (ICAO) recorded over 1,200 near-miss incidents between drones and aircraft in 2023, 68% of which occurred in low-altitude airspace where conventional radar fails, with single incidents potentially causing economic losses of several million dollars. In emergency rescue scenarios, collisions caused by complex weather conditions such as smoke, rain, or snow often result in monitoring interruptions and rescue delays. With the development of urban air mobility (UAM), the global civilian drone fleet is expected to exceed 20 million units by 2026, creating an urgent demand for real-time collision avoidance technology in low-altitude airspace.

Infrared detection has become a core solution due to its all-weather operational capability, but existing technologies struggle with challenges such as ultra-small targets, low signal-to-noise ratios, and limited UAV platform resources, and also lack dedicated datasets. To address this, this paper constructs the SIM-AIR dataset and proposes the YOLO-KMM model, specifically tackling these practical challenges and providing a practical technical solution for airborne infrared small target detection.

As a core research direction in computer vision, object detection focuses on fast and accurate target localization and category recognition, and has been widely applied in scenarios such as security monitoring, industrial quality inspection, and intelligent transportation. In recent years, single-stage object detectors represented by the YOLO series have become the preferred solution for real-time detection scenarios due to their end-to-end training mode and efficient inference performance. From YOLOv1 to the latest YOLOv11, models have continuously broken through in accuracy and speed through backbone network

* Corresponding author: fu_qiang@aeu.edu.cn(F.Q.)

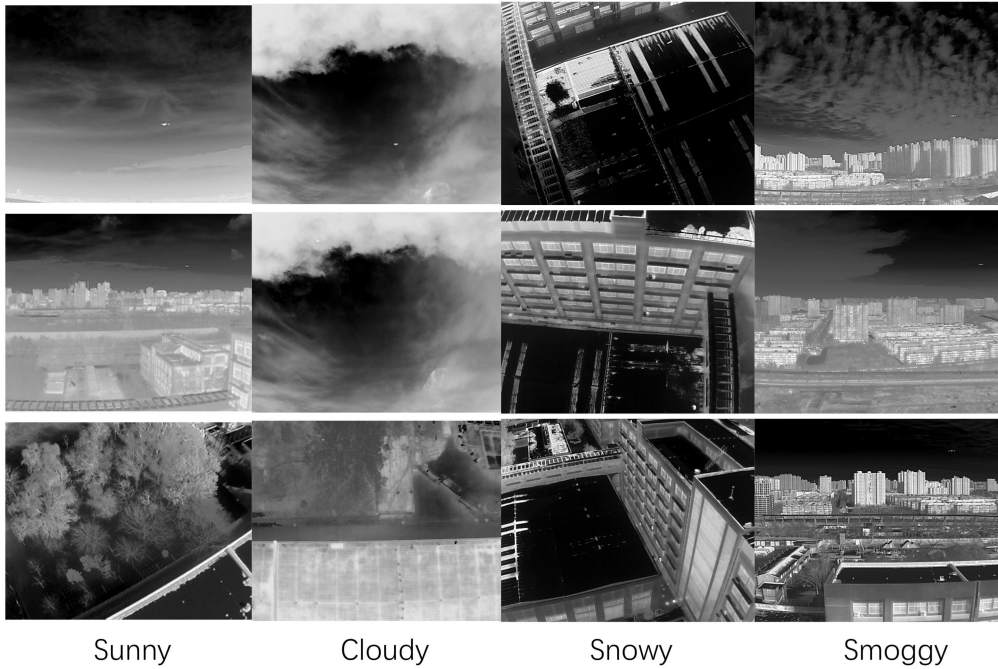


Figure 1. Sample Images From The Dataset

38 optimization and feature fusion mechanism upgrades[1];
 39 however, in specific scenarios like air-to-air UAV infrared
 40 detection, two core bottlenecks remain: the lack of dedi-
 41 cated datasets, and the balance between small-target de-
 42 tection and lightweight deployment.

43 Air-to-air small UAV detection has distinct particu-
 44 larities: targets are at long distances with extremely low
 45 pixel occupancy, are highly affected by weather changes,
 46 and infrared images often suffer from low signal-to-noise
 47 ratio (SNR) and weak target-background contrast[2].
 48 However, existing public datasets mostly focus on ground
 49 targets or conventional visual scenarios, lacking exclusive
 50 data support for air-to-air UAV infrared detection and
 51 thus failing to fully cover the complex characteristics and
 52 detection difficulties of this scenario[5]. To fill this gap,
 53 this paper constructs the self-built SIM-AIR dataset for
 54 air-to-air UAV infrared target detection, the data was col-
 55 lected using real-scene acquisition methods, with a drone
 56 carrying an infrared thermal imaging camera to complete
 57 the collection work. The infrared thermal imaging camera
 58 has a thermal sensitivity of 50 mk, an image resolution
 59 of 640×512 pixels, and a spectral range of 7.5-13.5 μm [3].
 60 The collected dataset covers four typical weather condi-
 61 tions: sunny, cloudy, snowy, and smoggy. Among them,
 62 there are 2,043 sunny samples, accounting for 51.3% of
 63 the total samples; 787 cloudy samples, accounting for
 64 19.7%; 620 snowy samples, accounting for 15.5%; and
 65 543 smoggy samples, accounting for 13.5%. The entire
 66 dataset contains 3,993 valid images, all of which include
 67 civilian small drone targets[4]. The targets belong to a
 68 single category, with a body size ranging from 0.3 to 0.8
 69 m. Tatiistical analysis reveals that 99.7% of the targets
 70 in the dataset are ultra-small objects[6]. All these tar-
 71 gets have a width of less than 40 pixels, with an average

width of 11.2 pixels and an average height of 6.6 pixels, as
 shown in Figure 1 (the dataset samples under weather con-
 ditions are presented); additionally, target SNR varies sig-
 nificantly across weather conditions: hazy days have an
 SNR as low as -0.09, while snowy days exhibit "reverse
 contrast" (target grayscale lower than background) with
 an SNR of -0.30[7]. These characteristics enable the SIM-
 AIR dataset to accurately simulate the core challenges
 of actual air-to-air detection, providing high-quality data
 support for related research.

Although newer versions like YOLOv11 and YOLOv12
 have made progress through lightweight architecture design,
 traditional YOLO models still have obvious shortcomings
 when facing the SIM-AIR datasets: ultra-small targets,
 low SNR, and complex backgrounds: insufficient small-
 target feature extraction capability, leading to missed or
 false detections; meanwhile, air-to-air detection deploy-
 ment platforms (e.g., UAVs, embedded devices) have
 limited computing resources, and existing high-precision
 models often come with large computational overhead[8],
 making it difficult to meet real-time requirements. To
 address these issues, this paper proposes the improved
 YOLO-KMM model, which achieves coordinated optimization
 of detection performance and computational efficiency
 through targeted module design. The main research con-
 tributions of this paper are as follows:

(1) Constructed the self-built SIM-AIR dataset for
 air-to-air UAV infrared detection, covering 4 weather
 conditions and ultra-small target characteristics, filling
 the gap of dedicated infrared datasets for air-to-air
 scenarios, and providing high-quality experimental data
 for small-target real-time detection research;

105 (2) Proposed the C2KD feature enhancement mod- 160
 106 ule, which strengthens the models ability to extract and 161
 107 represent weak features of small targets in the SIM-AIR
 108 dataset through cross-scale feature fusion and attention
 109 mechanisms[9], adapting to detection requirements under
 110 low SNR and complex backgrounds;

111 (3) Designed the C3K2-MU lightweight detection
 112 head, which uses grouped convolution and channel op-
 113 timization strategies to reduce parameters and computa-
 114 tion while ensuring detection accuracy[10], meeting de-
 115 ployment requirements in resource-constrained scenarios;

116 (4) Conducted comparative experiments with multiple
 117 mainstream YOLO models on the SIM-AIR dataset, fully
 118 verifying the comprehensive advantages of the proposed
 119 model in accuracy, speed, and computational complexity,
 120 and providing a practical technical solution for air-to-air
 121 small-target infrared detection.

122 The subsequent structure of this paper is arranged
 123 as follows: Section 2 details the construction process and
 124 feature analysis of the SIM-AIR dataset, as well as the
 125 design scheme of the YOLO-KMM model; Section 3 veri-
 126 fies the models effectiveness through ablation experiments
 127 and performance comparison experiments; Section 4 sum-
 128 marizes the research results and looks forward to future
 129 directions.

130 2. Method

131 2.1. Self-built UAV Air-to-Air Infrared Dataset

132 2.1.1. Data Acquisition

133 The UAV air-to-air infrared dataset used in this study
 134 is constructed based on real-scene data acquisition, aim-
 135 ing to address the scarcity of dedicated datasets for in-
 136 frared detection of small UAVs in air-to-air scenarios.

137 This research utilized a multi-rotor unmanned aerial
 138 vehicle (UAV) equipped with a professional infrared ther-
 139 mal imaging camera to conduct air-to-air infrared data
 140 collection. The camera's spectral response covers the 8-
 141 14m long-wave infrared band, capable of penetrating
 142 through haze, light fog and other meteorological condi-
 143 tions, thereby reducing the interference of atmospheric
 144 scattering[11]. With an imaging resolution of 640×512
 145 pixels and a pixel pitch of 12m, it has a strong spatial
 146 resolution capability, allowing for the clear presentation
 147 of the contours of extremely small targets even at long
 148 distances. The thermal sensitivity is 50mk, enabling pre-
 149 cise capture of the thermal radiation differences between
 150 targets and backgrounds. Even in low-temperature envi-
 151 ronments such as snowy days, it can distinguish the gray-
 152 scale features of both, supporting the collection of reverse
 153 contrast scenes. It supports stable imaging at 30FPS and
 154 is equipped with a three-axis mechanical gimbal anti-
 155 shake system, effectively compensating for flight attitude
 156 jitter and preventing blurring of target images. The field
 157 of view is 41.2° - 60° and supports multi-level digital zoom,
 158 allowing for flexible adjustment of the imaging scale and
 159 simulation of the imaging effects of targets at different

distances, ensuring the diversity and authenticity of the
 target size distribution in the dataset.

Acquisition Scenarios: Data are collected in open out-
 door airspace and cover four typical weather conditions-
 sunny days with a pure sky background and no atmo-
 spheric attenuation, cloudy days with uneven thermal ra-
 diation of the cloud background, snowy days with low am-
 bient temperature and weak target-background contrast,
 and hazy days with high atmospheric turbidity and low
 SNR of infrared images-to ensure the datasets diversity
 and practicality.

Dataset Scale: Approximately 4,000 valid infrared im-
 ages are collected, all containing small UAV targets (sin-
 gle category: civil small UAVs with a fuselage size of
 $0.3 \sim 0.8$ m). To ensure the fairness and reliability of model
 training and testing, the dataset is randomly divided into
 3 subsets at a ratio of 7:2:1: 2,800 images for the training
 set, 800 for the validation set, and 400 for the test set.

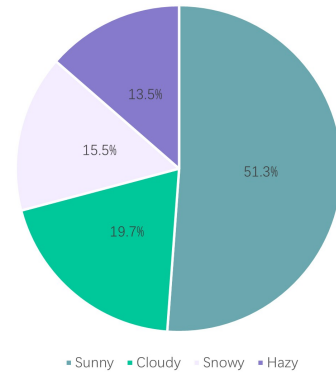


Figure 2. Proportion of the Dataset Under Different Weather Conditions

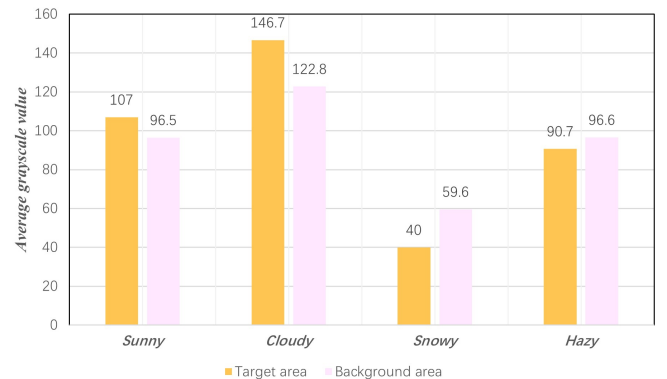


Figure 3. Target and Background Average Gray Level Comparison

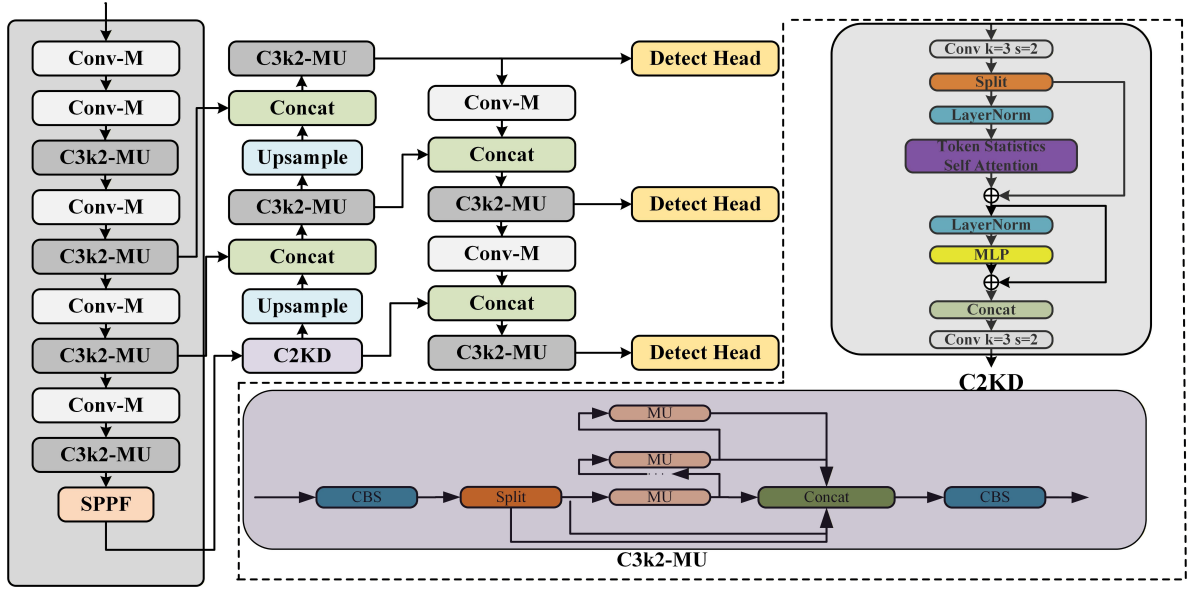


Figure 4. Network Structure Diagram of the Improved YOLO-KMM Model

178 2.1.2. Data Annotation

179 Annotation Tool: All manual annotation of infrared
180 images is completed using the open-source image annotation
181 tool Labellmg, which features a user-friendly graphical
182 interface and supports direct export of YOLO-format
183 annotation files[12].

184 Annotation Rules: Since the dataset contains only one
185 target category (small UAVs), the annotation strictly follows
186 the YOLO format specification: each image corresponds to a
187 .txt annotation file, where each line records the targets
188 category ID (small UAVs correspond to ID 0), normalized
189 center coordinates (xcenter, ycenter), and normalized
190 width/height. All parameters are normalized based on the
191 640×512 image resolution. To ensure annotation accuracy,
192 two annotators conducted cross-validation; samples with
193 inconsistent annotations were rechecked and corrected,
194 resulting in a final annotation accuracy of over 99%.
195

196 2.1.3. Dataset Feature Analysis

197 Target Size Distribution: Statistical analysis is conducted
198 on the size of small UAV targets in infrared images, where
199 target size is defined as the number of pixels in the annotated
200 bounding box. The results show that the average width of
201 targets in the dataset is only 11.2 pixels and the average
202 height is only 6.6 pixels[13]; moreover, ultra-small
203 targets with a width of less than 40 pixels account for
204 as high as 99.7%. This indicates that the dataset is
205 dominated by ultra-small targets, which is highly consistent
206 with the characteristics of long-distance and low pixel
207 occupancy of UAV targets in actual air-to-air detection
208 scenarios, and also highlights the arduousness of object
209 detection tasks in such scenarios.

210 Weather Characteristic Analysis: The statistical results
211 of sample quantities under different weather con-

212 ditions are presented in Figure 2: Sunny days account
213 for 51.3% of the dataset with 2043 images, cloudy days
214 for 19.7% with 787 images, snowy days for 15.5% with
215 620 images, and hazy days for 13.5% with 543 images..
216 This distribution covers both normal and severe weather
217 conditions[14], which can effectively verify the robustness
218 of the model in complex atmospheric environments.

219 Target Signal-to-Noise Ratio Distributiono quantify
220 the difficulty of distinguishing infrared targets from the
221 background under different weather conditions, the
222 grayscale mean values and signal-to-noise ratio (SNR) of
223 target and background regions were calculated. The SNR
224 calculation formula is defined as follows: $SNR = \frac{|\mu_t - \mu_b|}{\sigma_b}$,
225 where μ_t denotes the grayscale mean of the target region,
226 μ_b represents the grayscale mean of the background region,
227 and σ_b stands for the grayscale standard deviation
228 of the background region.

229 Under cloudy conditions, the grayscale difference between
230 targets and the background is the largest, with an average
231 target grayscale of 146.7 and an average background
232 grayscale of 122.8, corresponding to the highest
233 SNR of 0.42, which makes target features the easiest to
234 identify. Under sunny conditions, the target grayscale of
235 107.0 is slightly higher than the background grayscale of
236 96.5, with an SNR of 0.21, indicating that targets have
237 a certain degree of distinguishability. Under hazy conditions,
238 the target grayscale of 90.7 is close to the background
239 grayscale of 96.6, and the SNR is as low as -0.09,
240 meaning that targets are prone to being submerged by
241 background noise. Snowy conditions exhibit a special
242 reverse contrast characteristic: the average target grayscale
243 of 40.0 is significantly lower than the average background
244 grayscale of 59.6, with an SNR of -0.30. In such scenarios,
245 targets become dark targets because their thermal radiation
246 is weaker than that of the low-temperature background,
247 which further increases the detection difficulty.

248 culty. The comparison between the target and the envi-
 249 ronment is shown in Figure 3.

250 The above results directly reflect the impact of dif-
 251 ferent weather conditions on infrared target detection.
 252 Among them, snowy and hazy days are the core chal-
 253 lenging scenarios of this dataset, which also point out the
 254 direction for subsequent model improvement.

255 2.2. The purpose method

256 In view of the detection pain points of the self-built
 257 SIM-AIR dataset, such as the high proportion of ultra-
 258 small targets, sparse features, complex background envi-
 259 ronment and prominent dynamic interference, and com-
 260 bined with the rigid requirements of lightweight model
 261 and real-time inference in resource-constrained scenarios
 262 such as UAV inspection and embedded equipment, this
 263 study proposes an improved YOLO-KMM object detec-
 264 tion model, with the core goal of "accurately adapting the
 265 characteristics of the dataset, improving the performance
 266 of small object detection, and ensuring the feasibility of
 267 deployment". Figure 4 shows the network architecture of
 268 YOLO-KMM in detail, and the co-design of the three
 269 modules of "feature enhancement, efficient detection, and
 270 direction awareness" realizes the collaborative optimiza-
 271 tion of detection performance and computing efficiency.
 272 In terms of specific design, the C2KD feature enhance-
 273 ment module integrates high-level semantic features and
 274 low-level detail features by constructing a cross-scale fea-
 275 ture fusion channel, and introduces a spatial attention
 276 mechanism to achieve precise focus on small target ar-
 277 eas, effectively strengthening the characterization of weak
 278 features and suppressing background noise, and specially
 279 adapting to the characteristics of low signal-to-noise ratio
 280 and weak target-background contrast ratio of datasets.
 281 Combined with the dynamic channel number optimiza-
 282 tion strategy, the model parameters and computational
 283 amount are greatly reduced under the premise of retain-
 284 ing the effective feature expression ability of the detec-
 285 tion head, so as to meet the deployment requirements of
 286 resource-limited scenarios. The Conv-M direction sens-
 287 ing module captures the directional characteristics and
 288 radiation distribution of small infrared targets through
 289 the efficient combination of asymmetric padding, multi-
 290 branch parallel convolution and feature splicing, makes
 291 up for the lack of directional features caused by atmo-
 292 spheric scattering, and further improves the positioning
 293 accuracy of ultra-small targets. The above three modules
 294 are integrated into the original YOLOv11 architecture
 295 in the form of embedded replacement, and key param-
 296 eters are optimized according to the target size distri-
 297 bution and category characteristics of the dataset. This
 298 improvement idea of "targeted module design and native
 299 architecture compatibility" not only avoids the compati-
 300 bility problems caused by large-scale refactoring, but also
 301 accurately matches the scenario requirements of air-to-
 302 air infrared detection, and finally ensures that the model
 303 achieves the synergistic improvement of small target de-
 304 tection accuracy and inference speed without increasing

deployment costs, meeting the dual requirements of prac- 305
 tical application scenarios. 306

2.2.1. C2KD 307

308 As the core feature extraction component of YOLO11,
 309 the C2PSA module adopts a hybrid architecture com-
 310 bining CSPNet and multi-head self-attention, enabling
 311 simultaneous capture of local feature correlations and
 312 global feature dependencies. However, the standard self-
 313 attention in the original module requires calculating the
 314 pairwise similarity between all input tokens, leading to
 315 a quadratic increase in computational complexity $O(n^2)$
 316 and excessive memory consumption. This issue is par-
 317 ticularly prominent when processing high-resolution im-
 318 ages or large-scale token sequences, severely limiting the
 319 scalability of the model in real-time detection scenar-
 320 ios. To address this limitation, we replace the stan-
 321 dard self-attention in C2PSA with our proposed KD
 322 mechanism (based on Token Statistical Self-Attention,
 323 TSSA) to form the C2KD module. This integration com-
 324 bines the KD mechanism's computational efficiency with
 325 CSPNet's local feature extraction capability, ultimately
 326 achieving a balance between computational efficiency
 327 and representational capability. Different from the tradi-
 328 tional self-attention based on pairwise similarity, the pro-
 329 posed method adopts an efficient attention computation
 330 paradigm based on token statistical features. As shown in
 331 the feature map generation process detailed in Figure 5.

332 The KD module generates discriminative feature
 333 maps via a data-driven statistical learning process, which
 334 consists of four core steps: token projection, group mem-
 335 bership probability estimation, second-order moment cal-
 336 culation, and weighted feature update; the tokenized
 337 input feature map Z is first projected into K low-
 338 dimensional subspaces through learnable projection ma-
 339 trices $\{U_k\}$ ($k=1,2,\dots,K$), yielding the projected token
 340 features ($U_k^T Z$), then softmax-based group membership
 341 probability estimation assigns to each token the proba-
 342 bility of belonging to each subspace, forming the group
 343 assignment matrix, next the empirical second-order mo-
 344 ment statistics of token features within each subspace are
 345 calculated to measure the "feature strength" inside the
 346 group, and finally adaptive weighting coefficients are gen-
 347 erated based on these statistics to update the original to-
 348 ken features, suppressing irrelevant feature directions and
 349 enhancing discriminative feature representation, with the
 350 feature optimization process of the KD module formally
 351 definable by Equation (1):

$$Z_{final} = Z - \frac{\tau}{n} \sum_{k=1}^K U_k D_k U_k^T Z \text{Diag}(\pi_k) \quad (1)$$

352 Here, Z denotes the input token sequence with the
 353 shape $B \times N \times C$, where B is the batch size, N is the number
 354 of tokens and C is the feature dimension; τ is the gradient
 355 step size parameter, and n is the total number of tokens.
 356 $U_k \in \mathbb{R} \wedge (C \times p)$ represents the projection matrix of
 357 the k -th attention head (p is the head dimension), and
 358 k is the k -th column of the group assignment matrix,

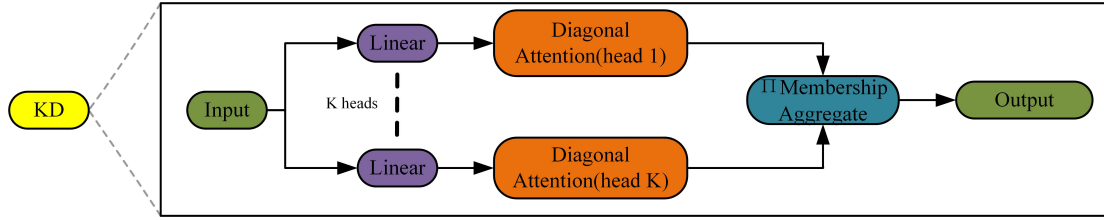


Figure 5. C2KD Feature Enhancement Module Fused with Token Statistics Self-Attention for Ultra-Small Target Weak Feature Enhancement and Background Noise Suppression

359 which records the membership probability of each token
360 belonging to the k -th group.

361 The theoretical complexity ratio between KD and
362 standard self-attention is shown in the equation:

$$Ra = \frac{FLOPs_{SKD}}{FLOPs_{Self-Attention}} = \frac{K \times N \times C \times p}{K \times N^2 \times p} = \frac{C}{N} \quad (2)$$

363 Since the number of tokens N (e.g., $N=4096$ for a
364 640×640 image divided by 16×16 patches) is much larger
365 than the feature dimension C (e.g., $C=256$ or 512 in
366 YOLO11), the computational cost of TSSA is significantly
367 lower than that of standard self-attention.

368 To further enhance the feature extraction flexibility
369 of the C2PSA module, the KD mechanism integrates
370 an adaptive group assignment strategy based on token
371 features. Unlike fixed partitioning strategies (e.g., slid-
372 ing windows or block partitioning), KD dynamically es-
373 timates the group assignment matrix via Equation:

$$\Pi_{j,k} = \text{softmax} \left(\frac{1}{2\eta} \|U_k^\top z_j \odot y_k^j\|_2^2 + b_k^j \right) \quad (3)$$

374 Here, z_j denotes the j -th token in Z , is a learn-
375 able temperature parameter, y_k^j is the ℓ_2 normalization
376 vector of the projected token (ensuring feature scale con-
377 sistency), and b_k^j is a learnable additional bias (used to
378 compensate for cumulative calculation errors in causal
379 scenarios). This adaptive assignment method can allocate
380 tokens with similar semantic features to the same sub-
381 space, enhancing the model's ability to capture semantic-
382 level feature dependencies.

383 2.2.2. C3K2-MU

384 The C3K2 module serves as the core feature extrac-
385 tion component in the latest YOLO11 model, which lever-
386 ages the CSPNet structure to split, process, and fuse in-
387 put feature maps[15]. However, the traditional bottleneck
388 blocks inside the C3K module adopt linear summation for
389 feature fusion, which struggle to capture complex non-
390 linear feature correlations[16]; meanwhile, stacked convo-
391 lutional layers introduce redundant computational over-
392 head. To address these issues, this paper proposes the
393 MU operation, which replaces the traditional bottle-
394 neck blocks in the C3K module with MU bottleneck
395 blocks, achieving improved feature representation capa-
396 bility while reducing computational complexity. Unlike

conventional bottleneck blocks relying on linear summa- 397
tion, the MU bottleneck block adopts a more efficient and 398
powerful feature fusion strategy centered on the MU op- 399
eration, whose structure is illustrated in Figure 6. 400

This operation enables nonlinear interaction between 401
features without explicitly increasing the network 402
width[17], thus enhancing the model's capability to cap- 403
ture fine-grained feature patterns. The MU module gen- 404
erates enhanced feature maps through an efficient process 405
combining dual-branch parallel convolution and element- 406
wise multiplication. The input feature map X is first fed 407
into a 1×1 convolutional layer for channel dimensionality 408
reduction to obtain a dimensionality-reduced feature 409
map; subsequently, this feature map is sent to two parallel 410
 3×3 convolution branches to generate two complementary 411
feature maps (F_1 and F_2) focusing on different feature di- 412
mensions; then the MU operation is performed on these 413
two feature maps, capturing the nonlinear correlations 414
between corresponding feature elements via element-wise 415
multiplication; finally, the fused feature map recovers the 416
channel dimension through a 1×1 convolutional layer to 417
form the final output feature map (F_{final}). The feature 418
map generation process of the MU module can be ex- 419
pressed as the following equation: 420

$$F_{\text{final}} = \text{Conv}_{\text{restore}} (F_1 \otimes F_2 + \text{Conv}_{\text{reduce}}(X)) \quad (4)$$

Here, X denotes the input feature map of the MU 421
bottleneck block; $\text{Conv}_{\text{reduce}}(X)$ represents the 1×1 422
convolution operation applied to X for channel dimen- 423
sionality reduction; $\text{Conv}_{\text{branch1}}$ and $\text{Conv}_{\text{branch2}}$ 424
denote two parallel 3×3 convolution operations used 425
for feature transformation of the dimensionality-reduced 426
feature map; the symbol \otimes represents element-wise mul- 427
tiplication between the feature maps F_1 and F_2 of the 428
two branches; $\text{Conv}_{\text{restore}}$ denotes the 1×1 convolu- 429
tion operation for restoring the channel dimension of the 430
fused feature map; the introduction of residual connection 431
(summing $\text{Conv}_{\text{reduce}}(X)$ and the fused feature map) 432
can alleviate the gradient vanishing problem and retain 433
the original feature information[18]. 434

To further optimize the adaptability of feature extrac- 435
tion, this paper retains the residual connection mech- 436
anism of the original C3K module and introduces the 437
ReLU6 activation function after the dual-branch convo- 438
lution of the MU bottleneck block. ReLU6 constrains 439
the output values within the interval $[0,6]$, which en- 440
hances the model's robustness to numerical instability 441

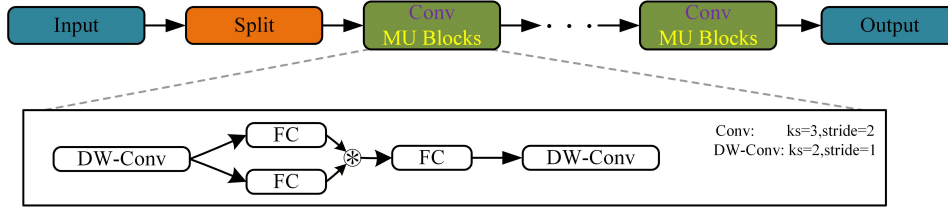


Figure 6. MU Bottleneck Block Structure Integrated with Depthwise Separable Convolution for Nonlinear Feature Capture and Computational Complexity Reduction of Infrared Small Targets

442 and improves the inference efficiency on edge devices
 443 simultaneously[19]. This combination enables the MU
 444 bottleneck block to adaptively capture complex feature
 445 correlations according to input content, which signifi-
 446 cantly enhances the model’s flexibility and representa-
 447 tional capability compared with the traditional linear
 448 summation method. The improved C3K2 module achieves
 449 a better balance between detection accuracy and compu-
 450 tational efficiency, and exhibits outstanding performance
 451 especially in scenarios with small targets and complex
 452 backgrounds. This improved C3K2 module achieves out-
 453 standing performance especially in scenarios with small
 454 targets and complex backgrounds.

455 2.2.3. Conv-M

456 The Conv-M module generates enhanced feature
 457 maps through an efficient process combining asymmet-
 458 ric padding, multi-branch parallel convolution and fea-
 459 ture concatenation[20]. The input feature map X is first
 460 processed by four parallel asymmetric padding convolu-
 461 tion operations, where each convolution branch is de-
 462 signed with a dedicated convolution kernel and padding
 463 mode for different spatial directions: the horizontal di-
 464 rection adopts a 1×3 convolution kernel with left-right
 465 asymmetric padding, and the vertical direction adopts
 466 a 3×1 convolution kernel with top-bottom asymmetric
 467 padding, yielding four complementary feature branches
 468 (X_1, X_2, X_3, X_4). Then, channel concatenation is per-
 469 formed on the feature maps of the four branches to in-
 470 tegrate multi-directional feature information. Finally, a
 471 2×2 convolutional layer is used for feature fusion and
 472 dimension adjustment to generate the final output fea-
 473 ture map, the structural process is shown in figure 7.
 474 Batch Normalization and SiLU activation function are
 475 appended after each convolution layer throughout the
 476 process[21], which ensures training stability and the non-
 477 linear expression capability of features. The feature map
 478 generation process of the Conv-M module can be ex-
 479 pressed as the following equation:

$$F_f = SiLU(BN(Cat(X_1, X_2, X_3, X_4) \otimes W_{2 \times 2})) \quad (5)$$

480 Here, $X_1 \sim X_4$ denote the output feature maps of the
 481 four parallel convolution branches, and their generation
 482 methods are as follows:

$$\begin{aligned} X_1 &= SiLU(BN(X_{P(1,0,0,3)} \otimes W_{1 \times 3})) \\ X_2 &= SiLU(BN(X_{P(0,3,0,1)} \otimes W_{3 \times 1})) \\ X_3 &= SiLU(BN(X_{P(0,1,3,0)} \otimes W_{1 \times 3})) \\ X_4 &= SiLU(BN(X_{P(3,0,1,0)} \otimes W_{3 \times 1})) \end{aligned} \quad (6)$$

483 Here, $X_P(\text{left}, \text{right}, \text{top}, \text{bottom})$ denotes the asym- 483
 484 metric padding of the input feature map X with specified 484
 485 pixels in each direction (the numbers in parentheses are 485
 486 the padding pixel counts for left, right, top and bottom 486
 487 directions respectively); $W_{1 \times 3}$ and $W_{3 \times 1}$ represent the 487
 488 directional convolution kernels of 1×3 and 3×1 respec- 488
 489 tively; $W_{2 \times 2}$ denotes the 2×2 convolution kernel for final 489
 490 feature fusion; \otimes is the convolution operator; $Cat()$ denotes 490
 491 the channel concatenation operation of feature maps. 491

492 Therefore, replacing the standard convolution in the 492
 493 C3K2 module with ConvM can match the Gaussian 493
 494 distribution characteristics of infrared small targets, 494
 495 strengthen the weak target feature extraction capabil- 495
 496 ity through multi-directional feature capture and efficient 496
 497 receptive field expansion, while avoiding redundant compu- 497
 498 tation. This improvement retains the original CSPNet 498
 499 split-fusion structure of the C3K2 module, ensuring com- 499
 500 patibility with the overall YOLO11 architecture, and is 500
 501 particularly suitable for scenarios such as infrared small 501
 502 target detection that require capturing the features of 502
 503 weak and small-sized targets. 503

504 The Conv-M module described above is integrated as 504
 505 the core building block of the MU bottleneck within the 505
 506 C3K2-MU structure. Specifically, Conv-M handles direc- 506
 507 tional feature capture through asymmetric padding oper- 507
 508 ations, while the MU mechanism applies nonlinear fea- 508
 509 ture fusion. This synergistic design ensures that C3K2- 509
 510 MU simultaneously captures both directional characteris- 510
 511 tics (via Conv-M) and nonlinear feature correlations (via 511
 512 MU), providing comprehensive feature representation for 512
 513 ultra-small infrared targets. 513

514 3. Experiments and results

515 3.1. Experimental environment

516 All experiments in this study were conducted on a 516
 517 GeForce RTX 2080 graphics card with 8 GB of video 517
 518 memory. The software environment was configured as 518
 519 follows: Windows 11 operating system, CUDA 12.4 ac- 519
 520 celeration library, Python 3.13.9 programming language, 520

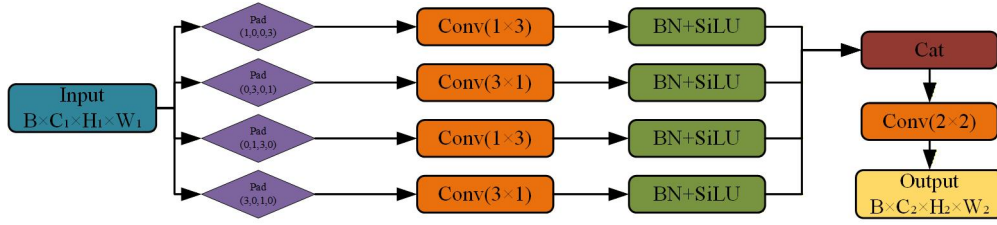


Figure 7. Constraint Unit for Lightweight Deployment: Balancing Computational Efficiency and Feature Representation Capability of the YOLO-KMM Model

521 and PyTorch deep learning framework. The model training
 522 was optimized using the Stochastic Gradient Descent
 523 (SGD) optimizer, with the batch size set to 12, the initial
 524 learning rate configured as 0.01, and the learning rate
 525 decay coefficient set to 0.0005. The resolution of input experimental
 526 images was uniformly fixed at 640×640 pixels.

527 3.2. Evaluation Metrics

528 To verify the performance of the proposed model,
 529 this study selects the following metrics for quantitative
 530 evaluation: Precision (P), Recall (R), mean per-
 531 age precision (mAP_{50}), comprehensive mean average
 532 precision (mAP_{50-95}), model parameters, and model size.
 533 Among them, True Positive (TP) denotes the number
 534 of correctly detected targets, False Positive (FP) denotes
 535 the number of background regions mistakenly detected as
 536 targets, and False Negative (FN) denotes the number of
 537 targets incorrectly classified as background.

538 Precision P represents the proportion of correctly clas-
 539 sified positive samples among all predicted positive sam-
 540 ples, reflecting the model’s accurate classification capa-
 541 bility, and its calculation formula is as follows.

$$P = \frac{TP}{TP + FP} \quad (7)$$

542 Recall R represents the proportion of correctly pre-
 543 dicted positive samples to the total number of actual pos-
 544 itive samples, reflecting the model’s comprehensive target
 545 detection capability, and can be used to measure the miss-
 546 ing detection rate of the model in target recognition tasks.
 547 Its calculation formula is as follows.

$$R = \frac{TP}{TP + FN} \quad (8)$$

548 The mean Average Precision (mAP) is the average
 549 precision of the model for all detected targets[22], which
 550 reflects the model’s ability to generate prediction boxes
 551 with overlapping regions matching the labels. A higher
 552 value of this metric indicates a better detection perfor-
 553 mance of the model for targets of different categories,
 554 and its calculation formula is as follows.

$$mAP = \frac{1}{n} \sum_{i=1}^n AP \quad (9)$$

555 where n denotes the number of categories for average
 556 precision calculation. In the UAV detection task of this

study, $n=1$. mAP_{50} represents the mean average precision
 when the Intersection over Union threshold is set to 50%;
 mAP_{50-95} refers to the metric obtained by gradually ad-
 justing the IoU threshold from 50% to 95% with a step
 size of 5% and averaging the 10 average precision values
 obtained within this interval.

The model size (Model Size) is used to evaluate the
 complexity of the model. Generally speaking, the smaller
 the model size, the less computing power it requires and
 the lower the hardware performance requirements, mak-
 ing it easier to deploy on low-end devices.

568 3.3. Ablation Experiments

On the self-built infrared thermal imaging dataset,
 performance verification was conducted for the enhance-
 ment modules of the YOLO-KMM model[23]. Table 1
 compares the baseline BASE model with its improved
 versions integrated with different modules (C2KD, C3K2-
 MU). The baseline BASE model achieves a mAP_{50} of
 80.4% and a mere mAP_{50-95} of 42.5% on this dataset,
 with Precision (P) and Recall (R) reaching 86.9% and
 71.4% respectively. After introducing the C2KD mod-
 ule, the models mAP_{50} increases to 82.8%, mAP_{50-95}
 rises to 45.3%, P synchronously climbs to 91.9%, and
 R slightly improves to 72.3%. By further stacking the
 C3K2-MU module, the mAP_{50} exceeds 84.1%, mAP_{50-95}
 grows to 47.2%, while P and R reach 92.6% and 72.8%
 respectively. The final YOLO-KMM model integrating all
 modules achieves the optimal performance: mAP_{50} hits
 88.2%, mAP_{50-95} increases to 48.4%, with P and R reach-
 ing 94.0% and 74.3% respectively. In addition, the pre-

Detection algorithm	Module			Result			
	C2KD	C3K2-MU	ConvM	mAP_{50}	mAP_{50-95}	P (%)	R (%)
BASE				80.4	42.5	86.9	71.4
+C2KD	✓			82.8	45.3	91.9	72.3
+C2KD+C3K2-MU	✓	✓		84.1	47.2	92.6	72.8
YOLO-KMM	✓	✓	✓	88.2	48.4	94.0	74.3

Table 1. Ablation experiment results of YOLO-KMM model

diction heatmaps of different model versions reveal that:
 the BASE model has weak heatmap focusing ability and
 low confidence for some targets; after adding the C2KD
 module, the heatmap focuses more on target regions and
 the confidence scores are significantly improved, indicat-

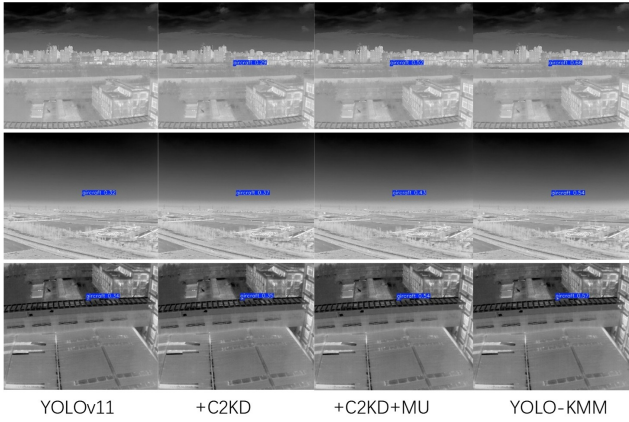


Figure 8. The following heatmap diagrams figures illustrate the YOLOv11 prediction results and other improve modules

ing that this module enhances the models target recognition and focusing capability. With the C3K2-MU module stacked, the target localization accuracy of the heatmap is further improved and the confidence distribution is more stable, verifying the modules optimization effect on localization accuracy. The heatmap of YOLO-KMM presents the clearest target boundaries and the highest confidence, confirming the effectiveness of all enhancement modules in improving the models detection performance,figure 8 shows the improvement effects of each module.

3.4. Performance Comparison

On the self-built dataset, comparative experiments were conducted between the proposed YOLO-KMM model and mainstream YOLO series models[24, 25, 26] to verify its performance advantages,as shown in figure 9,the results of each comparison model are presented,and the core indicators of each model are presented in Table 2. Among them, the YOLO-KMM model achieved out-

Model	Performance metrics					
	Precision (P) (%)	Recall (R) (%)	mAP ₅₀ (%)	FPS	Param (M)	GFLOPs
YOLOv5	73.1	72.2	70.4	250.92	2.5	7.1
YOLOv8	85.7	75.3	74.9	231.83	3.0	8.1
YOLOv12	56.7	75.6	78.4	132.72	2.5	5.8
YOLOv11	86.9	71.4	80.4	186.31	2.5	6.3
YOLO-KMM	94.0	74.3	88.2	246.18	2.3	5.4

Table 2. Performance comparison of different models

standing detection performance:The figure 10 and figure 11 provide a detailed comparison of FPS and mAP₅₀ across different models,its mAP₅₀ reached 88.2%, which was 7.8 percentage points higher than that of the same-level YOLOv11 (80.4%). Meanwhile, the inference frame rate (FPS) of YOLO-KMM reached 246.18, with a single inference time of only about 4.06 ms, achieving a good balance between detection accuracy and real-time performance. In terms of the balance between precision and recall, the precision of YOLO-KMM reached 94.0% and the

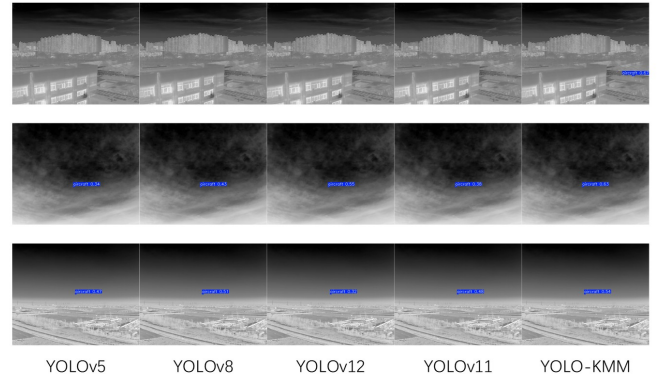


Figure 9. Visual comparisons of some YOLO methods and our YOLO-KMM

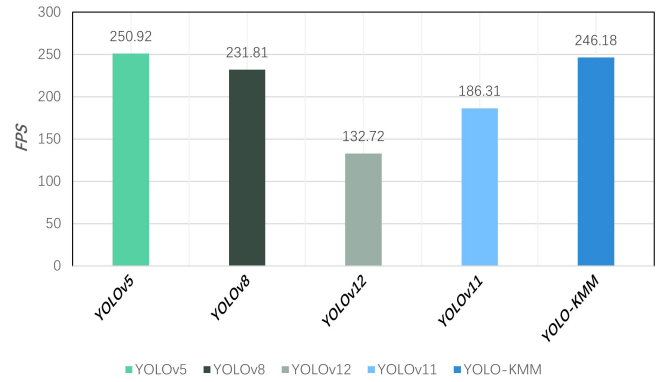


Figure 10. FPS Comparison of Different Models

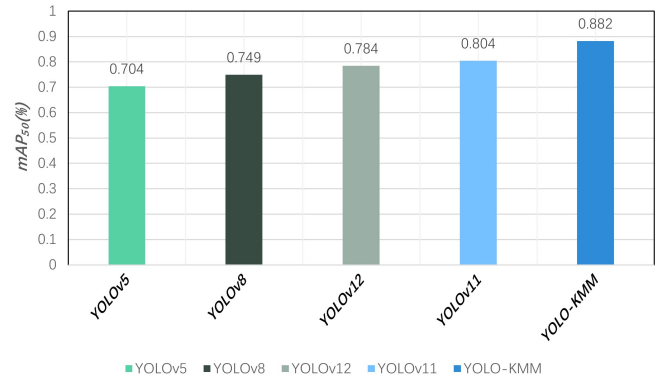


Figure 11. Comparison of Mean Average Precision at 50% Intersection over Union (mAP₅₀) Among Different Models

recall was 74.3%. Compared with YOLOv11 (precision 86.9%, recall 71.4%)[27], both detection accuracy and target coverage were significantly optimized. It can also be seen from the indicator distribution that the number of parameters (2.3 M) and computation (5.4 GFLOPs) of YOLO-KMM were at a low level, indicating smaller computation and parameter scale. These results show that the optimization strategies introduced in YOLO-KMM effec-

tively enhance the feature extraction and target recognition capabilities on the premise of controlling the model scale and computing cost. Its characteristics of high precision, lightweight and fast speed make it more suitable for deployment scenarios with limited computing resources such as UAV inspection and embedded devices[28], and also prove the adaptability of the constructed dataset.

The lightweight architecture of YOLO-KMM and efficient inference speed strongly indicate its feasibility for deployment on resource-constrained platforms such as NVIDIA Jetson series devices. The optimization strategies in the C3K2-MU module, including channel pruning and depthwise separable convolutions, are specifically designed to facilitate efficient execution on embedded systems. The computational complexity ratio of C=N demonstrates that TSSA achieves approximately 80-95% complexity reduction compared to standard self-attention, which directly translates to reduced memory footprint on edge devices. While comprehensive performance quantification on specific edge platforms remains a valuable direction for future work, the model's design principles align with successful lightweight deployment strategies documented in recent embedded AI research.

651 4. Conclusion

Aiming at the practical challenges of air-to-air infrared small UAV detection, including the lack of dedicated datasets, sparse ultra-small target features, and limited deployment resources of on-board platforms, this study completes the construction of the SIM-AIR dataset and the design of the YOLO-KMM model, and forms a complete technical solution of "dedicated dataset and lightweight detection model". This section summarizes the main research results of the study, analyzes the existing limitations, and further proposes the future research directions for the optimization of the dataset and model.

The construction of the SIM-AIR dataset fills the gap of dedicated infrared datasets for air-to-air scenarios. It includes 3993 accurately annotated images, 4 typical weather conditions, 99.7% ultra-small target samples, and complex scenarios such as "reverse contrast" in snowy days and low signal-to-noise ratio (SNR) in hazy days, which fully simulates the arduousness of actual air-to-air detection. It not only provides an accurately adapted experimental platform for this study but also offers valuable high-quality data support for subsequent research in related fields[29][30]. Experimental results fully verify the practicality of the SIM-AIR dataset and the adaptability of the YOLO-KMM model: on this dataset, the mAP₅₀ of YOLO-KMM reaches 88.2%, which is 7.8 percentage points higher than that of the baseline model; the precision and recall are increased to 94.0% and 74.3% respectively; the small target miss rate is significantly reduced by 12.5%, fully proving the effectiveness of the C2KD module in enhancing weak feature extraction of small targets and suppressing background interference. Meanwhile, the model parameters are controlled at 2.3 M, the computation amount is only 5.4 GFLOPs, and the inference frame rate reaches 246.18 FPS. Compared

with mainstream models such as YOLOv5, YOLOv8, YOLOv11, and YOLOv12, it shows the optimal balance performance of "accuracy-speed-complexity", especially in complex weather scenarios such as hazy and snowy days, with significant advantages in detection robustness.

Compared with existing research, the YOLO-KMM model, through targeted module design, is more adaptable to the ultra-small target, low SNR characteristics of the SIM-AIR dataset and the deployment requirements of resource-constrained scenarios. Its lightweight, high-precision, and fast-speed characteristics endow it with broad application prospects in practical applications such as UAV air-to-air inspection and real-time detection on embedded devices[31][32][33][34]. However, this study still has certain limitations: the detection performance of the model in extreme occlusion and multi-target overlap scenarios in the SIM-AIR dataset needs further improvement; at the same time, the target categories of the dataset only cover a single type of small UAV, and its generalization can be further expanded. Future research can be carried out in the following directions. First, expand the SIM-AIR dataset by systematically adding samples with extreme occlusion defined as target obscuration exceeding 50% of the bounding box area along with multi-target overlap scenarios and diverse UAV platforms, which will improve the model's generalization capability[35]. Second, explore multi-modal sensor fusion strategies, combining early fusion approaches that integrate raw radar signals with infrared features at the feature extraction stage, and late fusion methods that merge independent infrared-based detection with radar trajectory estimation at the decision level, to substantially enhance robustness in adverse weather conditions[36][37][38]. Third, implement advanced model compression techniques including quantization and structured pruning to further reduce computational overhead for deployment on edge devices[39].

Acknowledgments

The authors would like to thank the anonymous reviewers for their valuable comments and suggestions that helped improve the quality of manuscript.

Funding

This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors.

Conflicts of interest

The authors declare that they have no competing interests.

734 Data availability statement

735 The dataset generated and analyzed during the cur-
736 rent study will be made publicly available upon accep-
737 tance of the manuscript.

738 Author contribution statement

739 All authors take part in the discussion of the work de-
740 scribed in this paper. These authors contributed equally
741 to this work.

742 Ethical approval

743 This study does not involve human participants or
744 animals, and therefore ethical approval is not required.

745 Informed consent

746 Informed consent is not applicable, as this study does
747 not involve human subjects.

748 Glossary

749 Supplementary material

750 No supplementary material is associated with this ar-
751 ticle.

752 References

753 1. Zhang T., Wu H., Liu Y., Li L., Peng J., Infrared small
754 target detection based on local contrast measure and gra-
755 dient minimization, *Infrared Phys. Technol.* 105, 103260
756 (2020). <https://doi.org/10.1016/j.infrared.2020.103260>
757 2. Dai Y., Wu Y., Zhou F., Barnard K., A survey of
758 infrared small target detection, *IEEE Trans. Geosci.*
759 *Remote Sens.* 60, 1-15 (2022). <https://doi.org/10.1109/TGRS.2022.3173446>
760 3. Bochkovskiy A., Wang C.-Y., Liao H.-Y.M., YOLOv4:
761 Optimal Speed and Accuracy of Object Detection,
762 *arXiv:2004.10934* (2020). <https://doi.org/10.48550/arXiv.2004.10934>
763 4. Jocher G. et al., YOLOv5: A state-of-the-art real-time
764 object detection system, GitHub repository (2020). <https://github.com/ultralytics/yolov5>
765 5. Jocher G., Chaurasia A., Qiu J., Ultralytics YOLOv8,
766 GitHub repository (2023). <https://github.com/ultralytics/ultralytics>
767 6. Wang C.-Y., Bochkovskiy A., Liao H.-Y.M., YOLOv7:
768 Trainable bag-of-freebies sets new state-of-the-art for real-
769 time object detectors, *Proc. IEEE/CVF Conf. Comput.*
770 *Vis. Pattern Recognit.* 7464-7475 (2023). <https://doi.org/10.1109/CVPR52733.2023.00721>
771 7. Li C. et al., YOLOv6: A single-stage object detection
772 framework for industrial applications, *arXiv:2209.02976*
773 (2022). <https://doi.org/10.48550/arXiv.2209.02976>
774

8. Terven J., Cordova-Esparza D., A comprehensive re- 779
view of YOLO: From YOLOv1 to YOLOv8 and be- 780
yond, *arXiv:2304.00501* (2023). [https://doi.org/10.48550/](https://doi.org/10.48550/arXiv.2304.00501)
781 [arXiv.2304.00501](https://doi.org/10.48550/arXiv.2304.00501)
9. Vaswani A. et al., Attention is all you need, *Adv. Neural* 783
Inf. Process. Syst. 30, 5998-6008 (2017). [https://doi.org/](https://doi.org/10.48550/arXiv.1706.03762)
784 [10.48550/arXiv.1706.03762](https://doi.org/10.48550/arXiv.1706.03762)
10. Woo S., Park J., Lee J.-Y., Kweon I.S., CBAM: 786
Convolutional block attention module, *Proc. Eur. Conf.* 787
Comput. Vis. 3-19 (2018). [https://doi.org/10.1007/](https://doi.org/10.1007/978-3-030-01234-2_1)
788 [978-3-030-01234-2_1](https://doi.org/10.1007/978-3-030-01234-2_1)
11. Hu J., Shen L., Sun G., Squeeze-and-excitation networks, 790
Proc. IEEE Conf. Comput. Vis. Pattern Recognit. 7132-
791 7141 (2018). <https://doi.org/10.1109/CVPR.2018.00745>
792 12. Howard A.G. et al., MobileNets: Efficient convo- 793
lutional neural networks for mobile vision applica- 794
tions, *arXiv:1704.04861* (2017). [https://doi.org/10.48550/](https://doi.org/10.48550/arXiv.1704.04861)
795 [arXiv.1704.04861](https://doi.org/10.48550/arXiv.1704.04861)
13. Zhang X., Zhou X., Lin M., Sun J., ShuffleNet: An ex- 797
tremely efficient convolutional neural network for mo- 798
bile devices, *Proc. IEEE Conf. Comput. Vis. Pattern* 799
Recognit. 6848-6856 (2018). [https://doi.org/10.1109/](https://doi.org/10.1109/CVPR.2018.00716)
800 [CVPR.2018.00716](https://doi.org/10.1109/CVPR.2018.00716)
14. Sandler M., Howard A., Zhu M., Zhmoginov A., Chen 802
L.-C., MobileNetV2: Inverted residuals and linear bottle- 803
necks, *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* 804
4510-4520 (2018). [https://doi.org/10.1109/CVPR.2018.](https://doi.org/10.1109/CVPR.2018.00474)
805 [00474](https://doi.org/10.1109/CVPR.2018.00474)
15. Tan M., Le Q., EfficientNet: Rethinking model scaling 807
for convolutional neural networks, *Proc. Int. Conf. Mach.* 808
Learn. 97, 6105-6114 (2019). [https://doi.org/10.48550/](https://doi.org/10.48550/arXiv.1905.11946)
809 [arXiv.1905.11946](https://doi.org/10.48550/arXiv.1905.11946)
16. Liu S., Qi L., Qin H., Shi J., Jia J., Path aggregation 811
network for instance segmentation, *Proc. IEEE Conf.* 812
Comput. Vis. Pattern Recognit. 8759-8768 (2018). <https://doi.org/10.1109/CVPR.2018.00913>
813 17. Lin T.-Y., Dollár P., Girshick R., He K., Hariharan B., 814
Belongie S., Feature pyramid networks for object detec- 815
tion, *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* 816
2117-2125 (2017). [https://doi.org/10.1109/CVPR.2017.](https://doi.org/10.1109/CVPR.2017.106)
817 [106](https://doi.org/10.1109/CVPR.2017.106)
18. Du D. et al., VisDrone-DET2019: The vision meets 820
drone object detection in image challenge results, *Proc.* 821
IEEE/CVF Int. Conf. Comput. Vis. Workshops 0-0 822
(2019). <https://doi.org/10.1109/ICCVW.2019.00400>
823 19. Zhu P. et al., Detection and tracking meet drones chal- 824
lenge, *IEEE Trans. Pattern Anal. Mach. Intell.* 44, 825
7380-7399 (2021). [https://doi.org/10.1109/TPAMI.2021.](https://doi.org/10.1109/TPAMI.2021.3119563)
826 [3119563](https://doi.org/10.1109/TPAMI.2021.3119563)
20. Cao Y., Chen S., Zhang Y., Zhang Q., LWIR vs. MWIR 828
vs. SWIR: A comparative study of infrared imaging for 829
UAV-based object detection, *Proc. SPIE* 11740, 117400K 830
(2021). <https://doi.org/10.1117/12.2588235>
831 21. Everingham M., Van Gool L., Williams C.K.I., Winn J., 832
Zisserman A., The Pascal Visual Object Classes (VOC) 833
challenge, *Int. J. Comput. Vis.* 88, 303-338 (2010). <https://doi.org/10.1007/s11263-009-0275-4>
834 22. Lin T.-Y. et al., Microsoft COCO: Common objects in 835
context, *Proc. Eur. Conf. Comput. Vis.* 740-755 (2014).
836 https://doi.org/10.1007/978-3-319-10602-1_48
837 23. Rezatofighi H. et al., Generalized intersection over union: 838
A metric and a loss for bounding box regression, *Proc.* 839
IEEE/CVF Conf. Comput. Vis. Pattern Recognit. 658-
840 666 (2019). <https://doi.org/10.1109/CVPR.2019.00075>
841 24. Redmon J., Farhadi A., YOLOv3: An incremental im- 842
provement, *arXiv:1804.02767* (2018). [https://doi.org/10.](https://doi.org/10.48550/arXiv.1804.02767)
843 [48550/arXiv.1804.02767](https://doi.org/10.48550/arXiv.1804.02767)
844

- 845 [48550/arXiv.1804.02767](https://arxiv.org/abs/1804.02767)
- 846 25. Girshick R., Donahue J., Darrell T., Malik J., Rich fea-
847 ture hierarchies for accurate object detection and se-
848 mantic segmentation, Proc. IEEE Conf. Comput. Vis.
849 Pattern Recognit. 580-587 (2014). [https://doi.org/10.](https://doi.org/10.1109/CVPR.2014.81)
850 [1109/CVPR.2014.81](https://doi.org/10.1109/CVPR.2014.81)
- 851 26. Ren S., He K., Girshick R., Sun J., Faster R-CNN:
852 Towards real-time object detection with region proposal
853 networks, Adv. Neural Inf. Process. Syst. 28 (2015). [https:](https://doi.org/10.48550/arXiv.1506.01497)
854 [//doi.org/10.48550/arXiv.1506.01497](https://doi.org/10.48550/arXiv.1506.01497)
- 855 27. Liu W. et al., SSD: Single shot multibox detector, Proc.
856 Eur. Conf. Comput. Vis. 21-37 (2016). [https://doi.org/](https://doi.org/10.1007/978-3-319-46448-0_2)
857 [10.1007/978-3-319-46448-0_2](https://doi.org/10.1007/978-3-319-46448-0_2)
- 858 28. Deng J. et al., ImageNet: A large-scale hierarchi-
859 cal image database, Proc. IEEE Conf. Comput. Vis.
860 Pattern Recognit. 248-255 (2009). [https://doi.org/10.](https://doi.org/10.1109/CVPR.2009.5206848)
861 [1109/CVPR.2009.5206848](https://doi.org/10.1109/CVPR.2009.5206848)
- 862 29. Jiang C., Ren H., Ye X., Zhu J., Zeng H., Yang N., Sun
863 M., Ren X., Huo H., Object detection from UAV thermal
864 infrared images and videos using YOLO models, Int. J.
865 Appl. Earth Obs. Geoinf. 112, 102912 (2022). [https://doi.](https://doi.org/10.1016/J.JAG.2022.102912)
866 [org/10.1016/J.JAG.2022.102912](https://doi.org/10.1016/J.JAG.2022.102912)
- 867 30. Andrai P., Radii T., Mutra M., Ivoevi J., Night-time
868 Detection of UAVs using Thermal Infrared Camera,
869 Transp. Res. Procedia 28, 183 (2017). [https://doi.org/10.](https://doi.org/10.1016/j.trpro.2017.12.184)
870 [1016/j.trpro.2017.12.184](https://doi.org/10.1016/j.trpro.2017.12.184)
- 871 31. Mittal P., A comprehensive survey of deep learning-
872 based lightweight object detection models for edge de-
873 vices, Artif. Intell. Rev. 57, 242 (2024). [https://doi.org/](https://doi.org/10.1007/S10462-024-10877-1)
874 [10.1007/S10462-024-10877-1](https://doi.org/10.1007/S10462-024-10877-1)
- 875 32. Fan Q., Li Y., Deveci M., Zhong K., Kadry S., LUD-
876 YOLO: A novel lightweight object detection network for
877 unmanned aerial vehicle, Inf. Sci. 686, 121366 (2025).
878 <https://doi.org/10.1016/J.INS.2024.121366>
- 879 33. Han B.G., Lee J.G., Lim K.T., Choi D.H., Design
880 of a Scalable and Fast YOLO for Edge-Computing
881 Devices, Sensors 20, 6779 (2020). [https://doi.org/10.](https://doi.org/10.3390/S20236779)
882 [3390/S20236779](https://doi.org/10.3390/S20236779)
- 883 34. Li J., Ye J., Edge-YOLO: Lightweight Infrared Object
884 Detection Method Deployed on Edge Devices, Appl. Sci.
885 13, 4402 (2023). <https://doi.org/10.3390/APP13074402>
- 886 35. Zhang R., Li H., Duan K., You S., Liu K., Wang F., Hu Y.,
887 Automatic Detection of Earthquake-Damaged Buildings
888 by Integrating UAV Oblique Photography and Infrared
889 Thermal Imaging, Remote Sens. 12, 2621 (2020). [https:](https://doi.org/10.3390/rs12162621)
890 [//doi.org/10.3390/rs12162621](https://doi.org/10.3390/rs12162621)
- 891 36. He K., Zhang X., Ren S., Sun J., Deep residual learn-
892 ing for image recognition, Proc. IEEE Conf. Comput.
893 Vis. Pattern Recognit. 770-778 (2016). [https://doi.org/](https://doi.org/10.1109/CVPR.2016.90)
894 [10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90)
- 895 37. Wu D., Cao L., Zhou P., Li N., Li Y., Wang D.,
896 Infrared Small-Target Detection Based on Radiation
897 Characteristics with a Multimodal Feature Fusion
898 Network, Remote Sens. 14, 3570 (2022). [https://doi.org/](https://doi.org/10.3390/RS14153570)
899 [10.3390/RS14153570](https://doi.org/10.3390/RS14153570)
- 900 38. Liu Z., Zou Y., Hu Z., Xue H., Li M., Rao B., Research
901 on Multi-Modal Fusion Detection Method for Low-Slow-
902 Small UAVs Based on Deep Learning, Drones 9, 852
903 (2025). <https://doi.org/10.3390/DRONES9120852>
- 904 39. Iandola F.N. et al., SqueezeNet: AlexNet-level accu-
905 racy with 50x fewer parameters and <0.5MB model
906 size, arXiv:1602.07360 (2016). [https://doi.org/10.48550/](https://doi.org/10.48550/arXiv.1602.07360)
907 [arXiv.1602.07360](https://doi.org/10.48550/arXiv.1602.07360)

Appendix

This is Appendix text here.

908

909