

---

# HiddenObject: Modality-Agnostic Fusion for Multimodal Hidden Object Detection

---

**Anonymous Author(s)**

Affiliation

Address

email

## Abstract

Detecting hidden or partially concealed objects remains a fundamental challenge in multimodal environments, where factors like occlusion, camouflage, and lighting variations significantly hinder performance. Traditional RGB-based detection methods often *fail under such adverse conditions*, motivating the need for more robust, modality-agnostic approaches. In this work, we present **HiddenObject**, a novel fusion framework that integrates RGB, thermal, and depth data using a Mamba-based fusion mechanism. Our method *captures complementary signals across modalities*, enabling enhanced detection of obscured or camouflaged targets. Specifically, the proposed approach identifies modality-specific features and fuses them in a unified representation that generalizes well across challenging scenarios. We validate **HiddenObject** across multiple benchmark datasets, demonstrating *state-of-the-art or competitive performance* compared to existing methods. These results highlight the efficacy of our fusion design and expose key limitations in current unimodal and naïve fusion strategies. More broadly, our findings suggest that Mamba-based fusion architectures can significantly advance the field of multimodal object detection, especially under visually degraded or complex conditions.

## 1 Introduction

Rapid advancements in autonomous systems and robotic technologies have been motivated by the need

to enhance operational efficiency, mitigate labor shortages, and address complex real-world challenges

across diverse domains, such as industrial automation, agriculture, surveillance, search-and-rescue

operations, and security applications. As these scenarios grow increasingly sophisticated, accurately

detecting hidden or partially obscured objects becomes critical for effective decision-making and

operational success. Traditional object detection methods primarily depend on single-modality

imaging, notably RGB cameras, which are significantly limited by factors like lighting variations,

occlusions, camouflaged objects, and adverse environmental conditions [1].

Real-world objects often appear partially or fully concealed due to foliage, structural elements,

environmental clutter, or intentional concealment, posing substantial detection challenges for conven-

tional RGB imaging, which relies heavily on ideal lighting and unobstructed visibility. Conversely,

multimodal imaging techniques—such as thermal [2] and depth [3] imaging—offer robust solu-

tions by capturing complementary information even under challenging visual conditions. Thermal

imaging, for instance, detects infrared radiation independent of visible light, demonstrating strong

resilience against lighting variability and providing reliable performance in nighttime [4] or obscured

environments [5].

Integrating multiple imaging modalities thus presents significant advantages for handling complex

detection tasks, including partial occlusions [6], dense object clustering [7], and camouflage [8].

Objects exhibiting minimal color or texture contrast in RGB modalities often possess distinct signa-

tures in thermal, depth, or near-infrared (NIR) modalities, thereby greatly enhancing detectability. Specifically, thermal imaging differentiates objects based on thermal contrast, while depth imaging captures structural three-dimensional information, further facilitating the detection of hidden or obscured objects [9].

However, fusing multiple modalities into a cohesive detection framework introduces considerable technical challenges, such as modality misalignment, information redundancy, and computational efficiency. Existing approaches typically rely on static fusion strategies [10, 11] or emphasize single modalities [12], limiting their ability to interpret critical information effectively in complex or unfamiliar environments. To overcome these limitations, we introduce **HiddenObject**, a modality-agnostic detection framework utilizing a novel Mamba-based fusion mechanism coupled with a channel-aware decoder. Our approach efficiently extracts and seamlessly integrates critical information across diverse imaging modalities (e.g., RGB, thermal, depth), enabling robust and accurate detection of concealed objects in challenging scenarios.

As demonstrated in Figure 1, our framework maintains reliable detection performance even under severe occlusion and complex clutter conditions, a capability especially crucial in applications such as agriculture [13–15] and robotics [10], where target objects frequently remain hidden or obscured. The key contributions of our work are summarized as follows:

- We propose a Mamba-based fusion mechanism integrated with a channel-aware decoder to effectively extract and merge multimodal information from RGB, thermal, depth, and other modalities.
- Our detection pipeline robustly handles a wide range of hidden objects, including lightly obscured, partially occluded, heavily occluded, concealed, or camouflaged targets.
- Through extensive experimentation and analysis across multiple datasets, we validate the effectiveness and efficiency of our proposed framework, establishing a benchmark for future research into the capabilities of Mamba in multimodal learning.

In the following sections, we provide detailed descriptions of our approach, including the fusion mechanism, dataset selection, experimental protocols, and comprehensive performance evaluations, further emphasizing the practical applicability and versatility of our proposed **HiddenObject** framework.

## 2 Related Work

### 2.1 Object Detection and Segmentation Methods

In recent advancements, Carion et al. [16] pioneered the development of an end-to-end object detection framework utilizing transformers, termed DEtection TRansformer (DETR). This approach redefines object detection as a set prediction task, employing binary matching to generate one-to-one object predictions during training. Such a method streamlines the detection process by removing the necessity for manually crafted anchor boxes and non-maximum suppression (NMS) post-processing. Despite its benefits, DETR suffers from prolonged training convergence times, prompting the emergence of various enhanced DETR variants. Deformable DETR [17] improves training speed by predicting 2D anchor points and incorporating a deformable cross-attention mechanism to selectively sample features near reference points. Conditional DETR [18] separates content

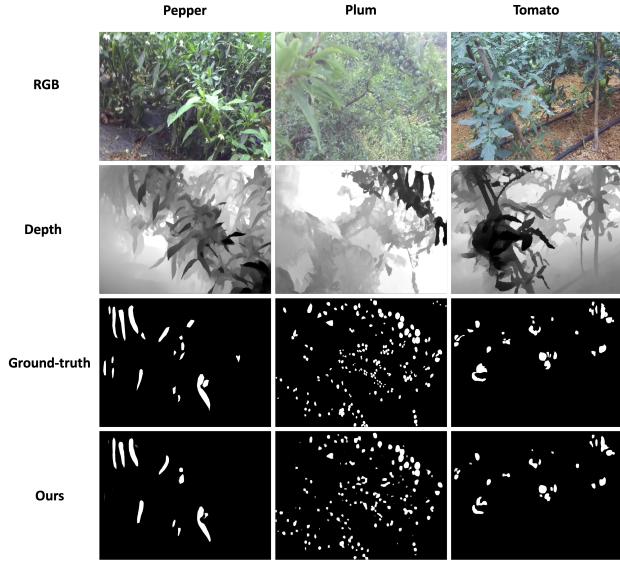


Figure 1: Our proposed method effectively addresses severe occlusion and densely distributed small objects, enhancing detection performance. The examples utilized are sourced from the ACOD-K12 dataset [6].

90 and positional operations, introducing conditional cross-attention to expedite convergence. Efficient  
91 DETR [19] enhances the pipeline’s efficiency by integrating dense detection with sparse set prediction.  
92 DAB-DETR [20] employs 4D reference points to refine bounding box predictions iteratively. DN-  
93 DETR [21] boosts training efficiency and label assignment by incorporating query denoising during  
94 training. DINO [22] consolidates these advancements into a robust DETR-based detection framework.  
95 For real-time performance, RT-DETR [23, 24] develops an end-to-end detector with an Efficient  
96 Hybrid Encoder and an IoU-aware Query Selection strategy.

97 Recent advancements in State Space Models (SSMs), particularly Mamba, offer efficient linear-time  
98 complexity for sequence modeling tasks, addressing limitations of Transformers in handling long  
99 sequences [25]. Extensions such as Vision Mamba (Vim) have adapted Mamba to visual tasks,  
100 improving speed and reducing memory consumption compared to traditional Vision Transformers  
101 like DeiT [26]. Further integration into object detection, including Mamba YOLO [27] and Fusion-  
102 Mamba [28], demonstrated significant gains in detection accuracy for multimodal data, crucial  
103 for identifying concealed objects. Multimodal fusion methods have also benefited from Mamba,  
104 particularly through Coupled Mamba, which enhances cross-modal representation consistency and  
105 inference speed, suggesting strong potential for applications in hidden object detection scenarios [29].

106 Earlier studies on infrared-visible object detection have largely relied on frameworks designed  
107 for single-modality object detection, typically categorized into one-stage detectors like Faster R-  
108 CNN [30] and two-stage detectors such as YOLO [31–33]. To integrate the complementary data  
109 from infrared and visible imagery, König et al. [34] proposed a fully convolutional fusion Region  
110 Proposal Network (RPN) that merges features from both modalities through concatenation, finding  
111 that mid-level fusion yields superior outcomes. Building on this, subsequent research developed  
112 CNN-based attention mechanisms to enhance the synergy between infrared and visible images  
113 [35–37]. Furthermore, some approaches incorporated transformer-based fusion modules to capture  
114 broader complementary relationships across these modalities [38–41]. Beyond direct feature fusion,  
115 certain methods utilized illumination data as global weights to combine infrared and visible features  
116 or to refine multi-branch detection outputs, mitigating the effects of noise [42, 43]. Recognizing that  
117 complementary traits vary by region, some studies employed semantic segmentation at the bounding  
118 box level [44, 45] or Region of Interest (ROI) prediction [46, 47] to direct regional fusion. Others  
119 leveraged confidence or uncertainty metrics to refine multi-branch predictions post-fusion [48, 49].  
120 To tackle modality misalignment, Zhang et al. [46] introduced the AR-CNN framework, aligning  
121 features across modalities using paired bounding box annotations to address object displacement.  
122 Similarly, Kim et al. [50] applied a multi-label learning strategy to adapt detection in misaligned  
123 scenarios. Recently, many multimodal methods have been introduced [51, 52]. Guo et al. [51]  
124 present a method for infrared-visible object detection that dynamically focuses on dominant modality  
125 objects and adaptively fuses complementary information, while Zhang et al. [52] proposed a novel  
126 end-to-end algorithm for multimodal fusion detection.

## 127 2.2 Datasets and Benchmarks

128 The KAIST Multispectral Pedestrian Dataset [53] consists of manually annotated RGB Color and  
129 Thermal imaging from the heat signatures of pedestrians from the perspective of a motor vehicle, with  
130 103,128 manual annotations. A multi-modality Neural Network-based fusion detection algorithm  
131 paper [54] collected RGB and Thermal data from a high-end FLIR T650sc infrared camera, capturing  
132 454 pairs of images containing construction defects from twenty years of wear. RGBT-1k [55]  
133 captured 1000 RGB-Thermal paired images from indoor and outdoor scenes with post-processed  
134 transformation steps for pixel-to-pixel alignment, which is an underlying concern for all RGB-X  
135 modality datasets. Osroosh et. al [56] developed an RGB-Thermal image collection technique, which  
136 in their methodology uniquely included a microclimate monitoring system with a rainfall simulator  
137 system to simulate different weather. The CitDet dataset [7] features 579 high-resolution images  
138 and 32,000 bounding box annotations of citrus trees, specifically targeting trees affected by the  
139 Huanglongbing bacterial infection disease, capable of significantly impacting citrus crop yield and  
140 quality through early detection and disease monitoring.

141 Thermal cameras, which capture infrared radiation emitted from object surfaces, have become increasingly  
142 popular due to their robustness under harsh conditions, widespread availability, and ability to  
143 provide unique information such as temperature. The proliferation of thermal cameras has led to the  
144 introduction of numerous thermal imaging datasets. The MultiSpectralMotion dataset [57] contains

145 indoor and outdoor thermal images collected through handheld devices along with ground-truth depth  
146 maps. The OdomBeyondVision dataset [58], in contrast, focuses on indoor environments using vari-  
147 ous platforms including handheld devices, unmanned ground vehicles (UGVs), and unmanned aerial  
148 systems (UAS). ViViD++ [59] provides outdoor thermal imagery captured using vehicle-mounted and  
149 handheld systems. The SubT-MRS dataset [60] comprises outdoor scenes collected from vehicles,  
150 UAS, legged robots, and handheld setups under varying degrees of visual degradation; it is most  
151 similar to our requirements but lacks imagery specific to wildland forests. Recent outdoor thermal  
152 stereo datasets include STheReO [61], MS2 [62], and FIReStereo [14]. MS2 offers thermal images  
153 captured in diverse environmental conditions, including visually challenging scenarios like rainy  
154 conditions. However, these datasets primarily focus on urban driving environments, inheriting limita-  
155 tions typical of driving datasets, such as constrained motion patterns and large baseline separation  
156 between cameras. FIReStereo [14] specifically targets small UAS applications aimed at developing  
157 depth estimation algorithms suitable for visually degraded environments.

## 158 References

- 159 [1] Zhengxia Zou, Keyan Chen, Zhenwei Shi, Yuhong Guo, and Jieping Ye. Object detection in 20  
160 years: A survey. *Proceedings of the IEEE*, 111(3):257–276, 2023.
- 161 [2] Mritunjay Rai, Tanmoy Maity, and RK Yadav. Thermal imaging system and its real time  
162 applications: a survey. *Journal of Engineering Technology*, 6(2):290–303, 2017.
- 163 [3] Alexandre Lopes, Roberto Souza, and Helio Pedrini. A survey on rgb-d datasets. *Computer*  
164 *Vision and Image Understanding*, 222:103489, 2022.
- 165 [4] Mario Bijelic, Tobias Gruber, Fahim Mannan, Florian Kraus, Werner Ritter, Klaus Dietmayer,  
166 and Felix Heide. Seeing through fog without seeing fog: Deep multimodal sensor fusion in  
167 unseen adverse weather. In *The IEEE Conference on Computer Vision and Pattern Recognition*  
(CVPR), June 2020.
- 168 [5] Hanzhe Teng, Yipeng Wang, Xiaoao Song, and Konstantinos Karydis. Multimodal dataset for  
localization, mapping and crop monitoring in citrus tree farms. In *Advances in Visual Computing*,  
169 pages 571–582, Cham, 2023. Springer Nature Switzerland. ISBN 978-3-031-47969-4.
- 170 [6] Liqiong Wang, Jinyu Yang, Yanfu Zhang, Fangyi Wang, and Feng Zheng. Depth-aware con-  
cealed crop detection in dense agricultural scenes. In *Proceedings of the IEEE/CVF Conference*  
171 *on Computer Vision and Pattern Recognition (CVPR)*, pages 17201–17211, June 2024.
- 172 [7] Jordan James, Heather Manching, Matthew Mattia, Kim Bowman, Amanda Hulse-Kemp, and  
173 William Beksi. Citdet: A benchmark dataset for citrus fruit detection. *IEEE Robotics and*  
174 *Automation Letters*, PP:1–8, 12 2024. doi: 10.1109/LRA.2024.3474473.
- 175 [8] Aditya Prakash, Kashyap Chitta, and Andreas Geiger. Multi-modal fusion transformer for  
end-to-end autonomous driving. In *Proceedings of the IEEE/CVF conference on computer*  
176 *vision and pattern recognition*, pages 7077–7087, 2021.
- 177 [9] Zhangyong Tang, Tianyang Xu, Xiao-Jun Wu, and Josef Kittler. Multi-level fusion for robust  
rgbt tracking via enhanced thermal representation. *ACM Transactions on Multimedia Computing,*  
178 *Communications and Applications*, 20(10):1–24, 2024.
- 179 [10] Qishen Ha, Kohei Watanabe, Takumi Karasawa, Yoshitaka Ushiku, and Tatsuya Harada. Mfnet:  
180 Towards real-time semantic segmentation for autonomous vehicles with multi-spectral scenes.  
181 In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages  
182 5108–5115, 2017. doi: 10.1109/IROS.2017.8206396.
- 183 [11] Yuxiang Sun, Weixun Zuo, and Ming Liu. Rtfnet: Rgb-thermal fusion network for semantic  
184 segmentation of urban scenes. *IEEE Robotics and Automation Letters*, 4(3):2576–2583, 2019.
- 185 [12] Shreyas S. Shivakumar, Neil Rodrigues, Alex Zhou, Ian D. Miller, Vijay Kumar, and Camillo J.  
186 Taylor. Pst900: Rgb-thermal calibration, dataset and segmentation network, 2019.

- 192 [13] Mahmoud Abdulsalam, Zakaria Chekakta, Nabil Aouf, and Maxwell Hogan. Fruity: A multi-  
 193 modal dataset for fruit recognition and 6D-Pose estimation in precision agriculture. In *2023*  
 194 *31st Mediterranean Conference on Control and Automation (MED)*, pages 144–149. IEEE, June  
 195 2023.
- 196 [14] Devansh Dhrafani, Yifei Liu, Andrew Jong, Ukcheol Shin, Yao He, Tyler Harp, Yaoyu Hu, Jean  
 197 Oh, and Sebastian Scherer. Firestereo: Forest infrared stereo dataset for uas depth perception  
 198 in visually degraded environments. *IEEE Robotics and Automation Letters*, 10(4):3302–3309,  
 199 2025. doi: 10.1109/LRA.2025.3536278.
- 200 [15] Connor Lee, Matthew Anderson, Nikhil Raganathan, Xingxing Zuo, Kevin Do, Georgia  
 201 Gkioxari, and Soon-Jo Chung. Cart: Caltech aerial rgb-thermal dataset in the wild. *arXiv*  
 202 preprint arXiv:2403.08997, 2024.
- 203 [16] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and  
 204 Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on*  
 205 *computer vision*, pages 213–229. Springer, 2020.
- 206 [17] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable  
 207 {detr}: Deformable transformers for end-to-end object detection. In *International Conference*  
 208 *on Learning Representations*, 2021.
- 209 [18] Depu Meng, Xiaokang Chen, Zejia Fan, Gang Zeng, Houqiang Li, Yuhui Yuan, Lei Sun, and  
 210 Jingdong Wang. Conditional detr for fast training convergence. In *Proceedings of the IEEE/CVF*  
 211 *International Conference on Computer Vision*, pages 3651–3660, 2021.
- 212 [19] Zhi Yao, Jiang Ai, Boxun Li, and Chi Zhang. Efficient detr: Improving end-to-end object  
 213 detector with dense prior. *arXiv preprint arXiv:2104.01318*, 2021.
- 214 [20] Shilong Liu, Feng Li, Hao Zhang, Xiao Yang, Xianbiao Qi, Hang Su, Jun Zhu, and Lei Zhang.  
 215 DAB-DETR: Dynamic anchor boxes are better queries for DETR. In *International Conference*  
 216 *on Learning Representations*, 2022.
- 217 [21] Feng Li, Hao Zhang, Shilong Liu, Jian Guo, Lionel M. Ni, and Lei Zhang. Dn-detr: Accelerate  
 218 detr training by introducing query denoising. In *Proceedings of the IEEE/CVF Conference on*  
 219 *Computer Vision and Pattern Recognition*, pages 13619–13627, 2022.
- 220 [22] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel Ni, and Heung-Yeung  
 221 Shum. DINO: DETR with improved denoising anchor boxes for end-to-end object detection. In  
 222 *The Eleventh International Conference on Learning Representations*, 2023.
- 223 [23] Yian Zhao, Wenyu Lv, Shangliang Xu, Jinman Wei, Guanzhong Wang, Qingqing Dang, Yi Liu,  
 224 and Jie Chen. Detrs beat yolos on real-time object detection. In *Proceedings of the IEEE/CVF*  
 225 *Conference on Computer Vision and Pattern Recognition*, pages 16965–16974, 2024.
- 226 [24] Wenyu Lv, Yian Zhao, Qinyao Chang, Kui Huang, Guanzhong Wang, and Yi Liu. Rt-detrv2:  
 227 Improved baseline with bag-of-freebies for real-time detection transformer, 2024. URL <https://arxiv.org/abs/2407.17140>.
- 228 [25] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. In  
 229 *First Conference on Language Modeling*, 2024.
- 230 [26] Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang.  
 231 Vision mamba: efficient visual representation learning with bidirectional state space model. In  
 232 *Proceedings of the 41st International Conference on Machine Learning*, ICML’24, 2024.
- 233 [27] Zeyu Wang, Chen Li, Huiying Xu, and Xinzhong Zhu. Mamba yolo: A simple baseline for  
 234 object detection with state space model. In *The 39th Annual AAAI Conference on Artificial*  
 235 *Intelligence*, 2025.
- 236 [28] Wenhao Dong, Haodong Zhu, Shaohui Lin, Xiaoyan Luo, Yunhang Shen, Xuhui Liu, Juan  
 237 Zhang, Guodong Guo, and Baochang Zhang. Fusion-mamba for cross-modality object detection.  
 238 *arXiv preprint arXiv:2402.16853*, 2024.

- 240 [29] Wenbing Li, Hang Zhou, Junqing Yu, Zikai Song, and Wei Yang. Coupled mamba: Enhanced  
 241 multimodal fusion with coupled state space model. In *The Thirty-eighth Annual Conference on*  
 242 *Neural Information Processing Systems*, 2024.
- 243 [30] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time  
 244 object detection with region proposal networks. *Advances in Neural Information Processing*  
 245 *Systems*, 28, 2015.
- 246 [31] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified,  
 247 real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and*  
 248 *Pattern Recognition*, pages 779–788, 2016.
- 249 [32] Joseph Redmon and Ali Farhadi. Yolo9000: Better, faster, stronger. In *Proceedings of the IEEE*  
 250 *Conference on Computer Vision and Pattern Recognition*, pages 7263–7271, 2017.
- 251 [33] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. Yolov7: Trainable bag-of-  
 252 freebies sets new state-of-the-art for real-time object detectors. In *Proceedings of the IEEE/CVF*  
 253 *Conference on Computer Vision and Pattern Recognition*, pages 7464–7475, 2023.
- 254 [34] Daniel König, Michael Adam, Christian Jarvers, Georg Layher, Heiko Neumann, and Michael  
 255 Teutsch. Fully convolutional region proposal networks for multispectral person detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*,  
 256 pages 49–56, 2017.
- 258 [35] Fang Qingyun and Wang Zhaokui. Cross-modality attentive feature fusion for object detection  
 259 in multispectral remote sensing imagery. *Pattern Recognition*, 130:108786, 2022.
- 260 [36] Kamil Roszyk, Michał R Nowicki, and Piotr Skrzypczyński. Adopting the yolov4 architecture  
 261 for low-latency multispectral pedestrian detection in autonomous driving. *Sensors*, 22(3):1082,  
 262 2022.
- 263 [37] Yue Cao, Junchi Bin, Jozsef Hamari, Erik Blasch, and Zheng Liu. Multimodal object detection  
 264 by channel switching and spatial attention. In *Proceedings of the IEEE/CVF Conference on*  
 265 *Computer Vision and Pattern Recognition*, pages 403–411, 2023.
- 266 [38] Fang Qingyun, Han Dapeng, and Wang Zhaokui. Cross-modality fusion transformer for  
 267 multispectral object detection. *arXiv preprint arXiv:2111.00273*, 2021.
- 268 [39] Haolong Fu, Shixun Wang, Puhong Duan, Changyan Xiao, Renwei Dian, Shutao Li, and  
 269 Zhiyong Li. Lraf-net: Long-range attention fusion network for visible-infrared object detection.  
 270 *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- 271 [40] Yaohui Zhu, Xiaoyu Sun, Miao Wang, and Hua Huang. Multi-modal feature pyramid transformer  
 272 for rgb-infrared object detection. *IEEE Transactions on Intelligent Transportation Systems*,  
 273 2023.
- 274 [41] Jifeng Shen, Yifei Chen, Yue Liu, Xin Zuo, Heng Fan, and Wankou Yang. Icafusion: Iterative  
 275 cross-attention guided feature fusion for multispectral object detection. *Pattern Recognition*,  
 276 145:109913, 2024.
- 277 [42] Kailai Zhou, Linsen Chen, and Xun Cao. Improving multispectral pedestrian detection by  
 278 addressing modality imbalance problems. In *Computer Vision–ECCV 2020: 16th European*  
 279 *Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVIII 16*, pages 787–803.  
 280 Springer, 2020.
- 281 [43] Xiaoxiao Yang, Yeqiang Qian, Huijie Zhu, Chunxiang Wang, and Ming Yang. Baanet: Learn-  
 282 ing bi-directional adaptive attention gates for multispectral pedestrian detection. In *2022*  
 283 *International Conference on Robotics and Automation (ICRA)*, pages 2920–2926. IEEE, 2022.
- 284 [44] Heng Zhang, Elisa Fromont, Sébastien Lefèvre, and Bruno Avignon. Multispectral fusion for  
 285 object detection with cyclic fuse-and-refine blocks. In *2020 IEEE International Conference on*  
 286 *Image Processing (ICIP)*, pages 276–280. IEEE, 2020.

- 287 [45] Heng Zhang, Elisa Fromont, Sébastien Lefèvre, and Bruno Avignon. Guided attentive fea-  
 288 ture fusion for multispectral pedestrian detection. In *Proceedings of the IEEE/CVF Winter*  
 289 *Conference on Applications of Computer Vision*, pages 72–80, 2021.
- 290 [46] Lu Zhang, Zhiyong Liu, Xiangyu Zhu, Zhan Song, Xu Yang, Zhen Lei, and Hong Qiao. Weakly  
 291 aligned feature fusion for multimodal object detection. *IEEE Transactions on Neural Networks*  
 292 and *Learning Systems*, 2021.
- 293 [47] Jung Uk Kim, Sungjune Park, and Yong Man Ro. Uncertainty-guided cross-modal learning for  
 294 robust multispectral pedestrian detection. *IEEE Transactions on Circuits and Systems for Video*  
 295 *Technology*, 32(3):1510–1523, 2021.
- 296 [48] Qing Li, Changqing Zhang, Qinghua Hu, Huazhu Fu, and Pengfei Zhu. Confidence-aware  
 297 fusion using dempster-shafer theory for multispectral pedestrian detection. *IEEE Transactions*  
 298 *on Multimedia*, 2022.
- 299 [49] Qing Li, Changqing Zhang, Qinghua Hu, Pengfei Zhu, Huazhu Fu, and Lei Chen. Stabilizing  
 300 multispectral pedestrian detection with evidential hybrid fusion. *IEEE Transactions on Circuits*  
 301 and *Systems for Video Technology*, 2023.
- 302 [50] Jiwon Kim, Hyeongjun Kim, Taejoo Kim, Namil Kim, and Yukyung Choi. Mlpd: Multi-label  
 303 pedestrian detector in multispectral domain. *IEEE Robotics and Automation Letters*, 6(4):  
 304 7846–7853, 2021.
- 305 [51] Junjie Guo, Chenqiang Gao, Fangcen Liu, Deyu Meng, and Xinbo Gao. Damsdet: Dynamic  
 306 adaptive multispectral detection transformer with competitive query selection and adaptive  
 307 feature fusion. In *European Conference on Computer Vision*, pages 464–481. Springer, 2024.
- 308 [52] Jiaqing Zhang, Mingxiang Cao, Weiyang Xie, Jie Lei, Daixun Li, Wenbo Huang, Yunsong  
 309 Li, and Xue Yang. E2E-MFD: Towards end-to-end synchronous multimodal fusion detection.  
 310 *Advances in Neural Information Processing Systems*, 37:52296–52322, 2024.
- 311 [53] Soonmin Hwang, Jaesik Park, Namil Kim, Yukyung Choi, and In So Kweon. Multispectral  
 312 pedestrian detection: Benchmark dataset and baselines. In *Proceedings of IEEE Conference on*  
 313 *Computer Vision and Pattern Recognition (CVPR)*, 2015.
- 314 [54] Xincong Yang, Runhao Guo, and Heng Li. Comparison of multimodal rgb-thermal fusion  
 315 techniques for exterior wall multi-defect detection. *Journal of Infrastructure Intelligence and*  
 316 *Resilience*, 2(2):100029, 2023. ISSN 2772-9915. doi: <https://doi.org/10.1016/j.iintel.2023.100029>.
- 318 [55] Santosh Kumar Panda and Pankaj Kumar Sa. Rgbt-1k: An rgb-thermal paired dataset, 2024.
- 319 [56] Yasin Osroosh and R. Troy Peters. Detecting fruit surface wetness using a custom-built low-  
 320 resolution thermal-rgb imager. *Computers and Electronics in Agriculture*, 157:509–517, 2019.  
 321 ISSN 0168-1699. doi: <https://doi.org/10.1016/j.compag.2019.01.023>.
- 322 [57] W. Dai, Y. Zhang, S. Chen, D. Sun, and D. Kong. A multi-spectral dataset for evaluating motion  
 323 estimation systems. In *2021 IEEE International Conference on Robotics and Automation*  
 324 (*ICRA*), pages 5560–5566. IEEE, 2021.
- 325 [58] P. Li, K. Cai, M. R. U. Saputra, Z. Dai, and C. X. Lu. Odombeyondvision: An indoor  
 326 multi-modal multi-platform odometry dataset beyond the visible spectrum. In *2022 IEEE/RSJ*  
 327 *International Conference on Intelligent Robots and Systems (IROS)*, pages 3845–3850. IEEE,  
 328 2022.
- 329 [59] A. J. Lee, Y. Cho, Y.-s. Shin, A. Kim, and H. Myung. Vivid++: Vision for visibility dataset.  
 330 *IEEE Robotics and Automation Letters*, 7(3):6282–6289, 2022.
- 331 [60] S. Zhao, Y. Gao, T. Wu, D. Singh, R. Jiang, H. Sun, M. Sarawata, Y. Qiu, W. Whittaker,  
 332 I. Higgins, Y. Du, S. Su, C. Xu, J. Keller, J. Karhade, L. Nogueira, S. Saha, J. Zhang, W. Wang,  
 333 C. Wang, and S. Scherer. Subt-mrs dataset: Pushing slam towards all-weather environments. In  
 334 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
 335 22647–22657, 2024.

336 [61] S. Yun, M. Jung, J. Kim, S. Jung, Y. Cho, M.-H. Jeon, G. Kim, and A. Kim. Sthereo: Stereo  
337 thermal dataset for research in odometry and mapping. In *2022 IEEE/RSJ International*  
338 *Conference on Intelligent Robots and Systems (IROS)*, pages 3857–3864. IEEE, 2022.

339 [62] U. Shin, J. Park, and I. S. Kweon. Deep depth estimation from thermal image. In *Proceedings*  
340 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages  
341 1043–1053, 2023.

342 **NeurIPS Paper Checklist**

343 **1. Claims**

344 Question: Do the main claims made in the abstract and introduction accurately reflect the paper's  
345 contributions and scope?

346 Answer: [Yes]

347 Justification: The paper proposes a new fusion mechanism, which is discussed in the Proposed  
348 Method and experimentation is performed and detailed in Section 4,

349 Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

359 **2. Limitations**

360 Question: Does the paper discuss the limitations of the work performed by the authors?

361 Answer: [Yes]

362 Justification: Limitations of the algorithm are covered in the conclusion.

363 Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

388 **3. Theory assumptions and proofs**

389 Question: For each theoretical result, does the paper provide the full set of assumptions and a  
390 complete (and correct) proof?

391 Answer: [NA]

392 Justification: The work produced in this paper is empirical, focusing on a new Mamba-based  
393 fusion architecture and experimentally evaluating the results.

- 394 Guidelines:
- 395 • The answer NA means that the paper does not include theoretical results.
- 396 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 397 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 398 • The proofs can either appear in the main paper or the supplemental material, but if they appear
- 399 in the supplemental material, the authors are encouraged to provide a short proof sketch to
- 400 provide intuition.
- 401 • Inversely, any informal proof provided in the core of the paper should be complemented by
- 402 formal proofs provided in appendix or supplemental material.
- 403 • Theorems and Lemmas that the proof relies upon should be properly referenced.

404 **4. Experimental result reproducibility**

405 Question: Does the paper fully disclose all the information needed to reproduce the main experi-  
406 mental results of the paper to the extent that it affects the main claims and/or conclusions of the  
407 paper (regardless of whether the code and data are provided or not)?

408 Answer: [Yes]

409 Justification: Formulas, along with the specific datasets, compute, and description of the Mamba  
410 architecture are provided.

411 Guidelines:

- 412 • The answer NA means that the paper does not include experiments.
- 413 • If the paper includes experiments, a No answer to this question will not be perceived well by the  
414 reviewers: Making the paper reproducible is important, regardless of whether the code and data  
415 are provided or not.
- 416 • If the contribution is a dataset and/or model, the authors should describe the steps taken to make  
417 their results reproducible or verifiable.
- 418 • Depending on the contribution, reproducibility can be accomplished in various ways. For  
419 example, if the contribution is a novel architecture, describing the architecture fully might  
420 suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary  
421 to either make it possible for others to replicate the model with the same dataset, or provide  
422 access to the model. In general, releasing code and data is often one good way to accomplish  
423 this, but reproducibility can also be provided via detailed instructions for how to replicate the  
424 results, access to a hosted model (e.g., in the case of a large language model), releasing of a  
425 model checkpoint, or other means that are appropriate to the research performed.
- 426 • While NeurIPS does not require releasing code, the conference does require all submissions  
427 to provide some reasonable avenue for reproducibility, which may depend on the nature of the  
428 contribution. For example
- 429 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to  
430 reproduce that algorithm.
- 431 (b) If the contribution is primarily a new model architecture, the paper should describe the  
432 architecture clearly and fully.
- 433 (c) If the contribution is a new model (e.g., a large language model), then there should either  
434 be a way to access this model for reproducing the results or a way to reproduce the model  
435 (e.g., with an open-source dataset or instructions for how to construct the dataset).
- 436 (d) We recognize that reproducibility may be tricky in some cases, in which case authors are  
437 welcome to describe the particular way they provide for reproducibility. In the case of  
438 closed-source models, it may be that access to the model is limited in some way (e.g.,  
439 to registered users), but it should be possible for other researchers to have some path to  
440 reproducing or verifying the results.

441 **5. Open access to data and code**

442 Question: Does the paper provide open access to the data and code, with sufficient instructions to  
443 faithfully reproduce the main experimental results, as described in supplemental material?

444 Answer: [No]

445 Justification: There is no provided open access to the data and code for now, but it is in the future  
446 intent to open source the code within the next calendar year.

447 Guidelines:

- 448     • The answer NA means that paper does not include experiments requiring code.  
 449     • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.  
 450  
 451     • While we encourage the release of code and data, we understand that this might not be possible,  
 452       so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless  
 453       this is central to the contribution (e.g., for a new open-source benchmark).  
 454     • The instructions should contain the exact command and environment needed to run to reproduce  
 455       the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.  
 456  
 457     • The authors should provide instructions on data access and preparation, including how to access  
 458       the raw data, preprocessed data, intermediate data, and generated data, etc.  
 459  
 460     • The authors should provide scripts to reproduce all experimental results for the new proposed  
 461       method and baselines. If only a subset of experiments are reproducible, they should state which  
 462       ones are omitted from the script and why.  
 463  
 464     • At submission time, to preserve anonymity, the authors should release anonymized versions (if  
 465       applicable).  
 466     • Providing as much information as possible in supplemental material (appended to the paper) is  
 467       recommended, but including URLs to data and code is permitted.

466     6. **Experimental setting/details**

467     Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters,  
 468       how they were chosen, type of optimizer, etc.) necessary to understand the results?

469     Answer: [Yes]

470     Justification: Details on a training and test details are specified in the experimentation.

471     Guidelines:

- 472     • The answer NA means that the paper does not include experiments.  
 473     • The experimental setting should be presented in the core of the paper to a level of detail that is  
 474       necessary to appreciate the results and make sense of them.  
 475     • The full details can be provided either with the code, in appendix, or as supplemental material.

476     7. **Experiment statistical significance**

477     Question: Does the paper report error bars suitably and correctly defined or other appropriate  
 478       information about the statistical significance of the experiments?

479     Answer: [No]

480     Justification: Our paper does not report error bars.

481     Guidelines:

- 482     • The answer NA means that the paper does not include experiments.  
 483     • The authors should answer "Yes" if the results are accompanied by error bars, confidence  
 484       intervals, or statistical significance tests, at least for the experiments that support the main claims  
 485       of the paper.  
 486     • The factors of variability that the error bars are capturing should be clearly stated (for example,  
 487       train/test split, initialization, random drawing of some parameter, or overall run with given  
 488       experimental conditions).  
 489     • The method for calculating the error bars should be explained (closed form formula, call to a  
 490       library function, bootstrap, etc.).  
 491     • The assumptions made should be given (e.g., Normally distributed errors).  
 492     • It should be clear whether the error bar is the standard deviation or the standard error of the  
 493       mean.  
 494     • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably  
 495       report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of  
 496       errors is not verified.  
 497     • For asymmetric distributions, the authors should be careful not to show in tables or figures  
 498       symmetric error bars that would yield results that are out of range (e.g. negative error rates).  
 499     • If error bars are reported in tables or plots, The authors should explain in the text how they were  
 500       calculated and reference the corresponding figures or tables in the text.

501 8. **Experiments compute resources**

502 Question: For each experiment, does the paper provide sufficient information on the computer  
503 resources (type of compute workers, memory, time of execution) needed to reproduce the experi-  
504 ments?

505 Answer: [Yes]

506 Justification: The paper includes information on the bare metal that the Mamba architecture, along  
507 with the experiments that are performed upon.

508 Guidelines:

- 509 • The answer NA means that the paper does not include experiments.
- 510 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud  
511 provider, including relevant memory and storage.
- 512 • The paper should provide the amount of compute required for each of the individual experimental  
513 runs as well as estimate the total compute.
- 514 • The paper should disclose whether the full research project required more compute than the  
515 experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it  
516 into the paper).

517 9. **Code of ethics**

518 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS  
519 Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

520 Answer: [Yes]

521 Justification: We have reviewed the NeurIPS Code of Ethics and checked that the paper conforms  
522 with the NeurIPS Code of Ethics.

523 Guidelines:

- 524 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 525 • If the authors answer No, they should explain the special circumstances that require a deviation  
526 from the Code of Ethics.
- 527 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration due  
528 to laws or regulations in their jurisdiction).

529 10. **Broader impacts**

530 Question: Does the paper discuss both potential positive societal impacts and negative societal  
531 impacts of the work performed?

532 Answer: [Yes]

533 Justification: Yes, the Broader Impacts section discusses the positive and negative societal impacts  
534 of the work performed.

535 Guidelines:

- 536 • The answer NA means that there is no societal impact of the work performed.
- 537 • If the authors answer NA or No, they should explain why their work has no societal impact or  
538 why the paper does not address societal impact.
- 539 • Examples of negative societal impacts include potential malicious or unintended uses (e.g.,  
540 disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deploy-  
541 ment of technologies that could make decisions that unfairly impact specific groups), privacy  
542 considerations, and security considerations.
- 543 • The conference expects that many papers will be foundational research and not tied to par-  
544 ticular applications, let alone deployments. However, if there is a direct path to any negative  
545 applications, the authors should point it out. For example, it is legitimate to point out that  
546 an improvement in the quality of generative models could be used to generate deepfakes for  
547 disinformation. On the other hand, it is not needed to point out that a generic algorithm for  
548 optimizing neural networks could enable people to train models that generate Deepfakes faster.
- 549 • The authors should consider possible harms that could arise when the technology is being used  
550 as intended and functioning correctly, harms that could arise when the technology is being used  
551 as intended but gives incorrect results, and harms following from (intentional or unintentional)  
552 misuse of the technology.

- 553     • If there are negative societal impacts, the authors could also discuss possible mitigation strategies  
554       (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for  
555       monitoring misuse, mechanisms to monitor how a system learns from feedback over time,  
556       improving the efficiency and accessibility of ML).

557   **11. Safeguards**

558   Question: Does the paper describe safeguards that have been put in place for responsible release of  
559   data or models that have a high risk for misuse (e.g., pretrained language models, image generators,  
560   or scraped datasets)?

561   Answer: [NA]

562   Justification: The paper is not a pretrained language model, image generator, scraped dataset, or  
563   anything that may have high risk for misuse.

564   Guidelines:

- 565     • The answer NA means that the paper poses no such risks.
- 566     • Released models that have a high risk for misuse or dual-use should be released with necessary  
567       safeguards to allow for controlled use of the model, for example by requiring that users adhere  
568       to usage guidelines or restrictions to access the model or implementing safety filters.
- 569     • Datasets that have been scraped from the Internet could pose safety risks. The authors should  
570       describe how they avoided releasing unsafe images.
- 571     • We recognize that providing effective safeguards is challenging, and many papers do not require  
572       this, but we encourage authors to take this into account and make a best faith effort.

573   **12. Licenses for existing assets**

574   Question: Are the creators or original owners of assets (e.g., code, data, models), used in the  
575   paper, properly credited and are the license and terms of use explicitly mentioned and properly  
576   respected?

577   Answer: [Yes]

578   Justification: We used published datasets that were correctly referenced and followed their terms  
579   of use in accordance with usage and citations.

580   Guidelines:

- 581     • The answer NA means that the paper does not use existing assets.
- 582     • The authors should cite the original paper that produced the code package or dataset.
- 583     • The authors should state which version of the asset is used and, if possible, include a URL.
- 584     • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 585     • For scraped data from a particular source (e.g., website), the copyright and terms of service of  
586       that source should be provided.
- 587     • If assets are released, the license, copyright information, and terms of use in the package should  
588       be provided. For popular datasets, [paperswithcode.com/datasets](http://paperswithcode.com/datasets) has curated licenses for  
589       some datasets. Their licensing guide can help determine the license of a dataset.
- 590     • For existing datasets that are re-packaged, both the original license and the license of the derived  
591       asset (if it has changed) should be provided.
- 592     • If this information is not available online, the authors are encouraged to reach out to the asset's  
593       creators.

594   **13. New assets**

595   Question: Are new assets introduced in the paper well documented and is the documentation  
596   provided alongside the assets?

597   Answer: [No]

598   Justification: The paper evaluates on several public datasets, and describes a framework of the  
599   approach.

600   Guidelines:

- 601     • The answer NA means that the paper does not release new assets.
- 602     • Researchers should communicate the details of the dataset/code/model as part of their sub-  
603       missions via structured templates. This includes details about training, license, limitations,  
604       etc.

- 605 • The paper should discuss whether and how consent was obtained from people whose asset is  
606 used.  
607 • At submission time, remember to anonymize your assets (if applicable). You can either create  
608 an anonymized URL or include an anonymized zip file.

609 **14. Crowdsourcing and research with human subjects**

610 Question: For crowdsourcing experiments and research with human subjects, does the paper  
611 include the full text of instructions given to participants and screenshots, if applicable, as well as  
612 details about compensation (if any)?

613 Answer: [NA]

614 Justification: The paper does not involve crowdsourcing experiments and research with human  
615 subjects.

616 Guidelines:

- 617 • The answer NA means that the paper does not involve crowdsourcing nor research with human  
618 subjects.  
619 • Including this information in the supplemental material is fine, but if the main contribution of  
620 the paper involves human subjects, then as much detail as possible should be included in the  
621 main paper.  
622 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other  
623 labor should be paid at least the minimum wage in the country of the data collector.

624 **15. Institutional review board (IRB) approvals or equivalent for research with human subjects**

625 Question: Does the paper describe potential risks incurred by study participants, whether such  
626 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals  
627 (or an equivalent approval/review based on the requirements of your country or institution) were  
628 obtained?

629 Answer: [NA]

630 Justification: The paper does not involve crowdsourcing experiments and research with human  
631 subjects.

632 Guidelines:

- 633 • The answer NA means that the paper does not involve crowdsourcing nor research with human  
634 subjects.  
635 • Depending on the country in which research is conducted, IRB approval (or equivalent) may be  
636 required for any human subjects research. If you obtained IRB approval, you should clearly  
637 state this in the paper.  
638 • We recognize that the procedures for this may vary significantly between institutions and  
639 locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for  
640 their institution.  
641 • For initial submissions, do not include any information that would break anonymity (if applica-  
642 ble), such as the institution conducting the review.

643 **16. Declaration of LLM usage**

644 Question: Does the paper describe the usage of LLMs if it is an important, original, or non-  
645 standard component of the core methods in this research? Note that if the LLM is used only for  
646 writing, editing, or formatting purposes and does not impact the core methodology, scientific  
647 rigorosity, or originality of the research, declaration is not required.

648 Answer: [NA]

649 Justification: Core research methods do not involve usage of LLM's.

650 Guidelines:

- 651 • The answer NA means that the core method development in this research does not involve LLMs  
652 as any important, original, or non-standard components.  
653 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what  
654 should or should not be described.