

PERCEPTUAL PIERCING: HUMAN VISUAL CUE-BASED OBJECT DETECTION IN LOW VISIBILITY CONDITIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

This study proposes a novel deep learning framework inspired by atmospheric scattering and human visual cortex mechanisms to enhance object detection under poor visibility scenarios such as fog, smoke, and haze. These conditions pose significant challenges for object recognition, impacting various sectors, including autonomous driving, aviation management, and security systems. The objective is to enhance the precision and reliability of detection systems under adverse environmental conditions. The research investigates the integration of human-like visual cues, particularly focusing on selective attention and environmental adaptability, to ascertain their impact on object detection's computational efficiency and accuracy. This paper proposes a multi-tiered strategy that integrates an initial quick detection process, followed by targeted region-specific dehazing, and concludes with an in-depth detection phase. The approach is validated using the Foggy Cityscapes, RESIDE- β (OTS and RTTS) datasets and is anticipated to set new performance standards in detection accuracy while significantly optimizing computational efficiency. The findings offer a viable solution for enhancing object detection in poor visibility and contribute to the broader understanding of integrating human visual principles into deep learning algorithms for intricate visual recognition challenges. The code for perceptual piercing is available here.

1 INTRODUCTION

Low-visibility conditions such as rain, snow, fog, smoke, or haze present significant challenges in various fields of computer vision and deep learning, such as autonomous vehicles, security and surveillance, maritime navigation, and agricultural robotics. The objective is to develop a deep-learning framework capable of recognizing objects using human visual cues under adverse visibility conditions. The motivation behind this project lies in addressing the substantial difficulties of identifying objects in low-visibility environments, a critical factor in enhancing airport operations during adverse weather.

Poor visibility often leads to aircraft delays, as planes face challenges in taxiing to their gates without clear visual guidance. This situation necessitates more ground support personnel to assist planes in docking, but due to limited ground staff availability, a bottleneck can occur, impeding the handling of multiple aircraft and resulting in further delays. These delays can escalate, potentially leading to flight cancellations. Although the initial motivation for this project is rooted in reducing delays in airport operations, the scope of the proposed machine learning model extends beyond airport scenarios to include a broad range of low-visibility environments. The following methods have been proposed:

- **Selective Region Enhancement:** Unlike uniform dehazing, focusing on specific regions can reduce processing time and prevent image quality degradation in areas where clarity might introduce false positives or where detail is not essential for current detection goals.
- **Integration with Object Detection:** By bridging the gap between image enhancement and object detection, we offer a cohesive approach that leverages the strengths of both methodologies, addressing the limitations of traditional, separate systems.

The above contributions are inspired by mechanisms of the human visual system, including selective attention, foveal and peripheral vision, human-eye adjustments to environmental conditions,

054 eye-tracking concepts, bottom-up signals based on sensory input, and top-down processes guided
055 by priors and current goals.

056 The rest of the article is organized as follows: Section 2 reviews the related work, outlining previous
057 studies and developments pertinent to object detection in low-visibility conditions and the integra-
058 tion of human visual cues into machine learning models. This section also highlights the gaps in
059 current research that our study aims to address. Section 3 describes the methodology of our study,
060 detailing the proposed deep-learning framework inspired by human visual signals, the selection cri-
061 teria for our datasets, and the experimental setup used to evaluate the model’s performance under
062 various low-visibility scenarios.

064 2 RELATED WORK

066 The field of navigation and detection in low-visibility conditions has seen significant advancements
067 through various methodologies including sensor fusion, visual cue integration, and computational
068 techniques. Aircraft landing has been a focus area, with studies exploring sensor fusion of visible
069 and virtual imagery (Liu et al., 2014) and visual-inertial navigation algorithms relying on synthetic
070 and real runway features (Zhang et al., 2018). For GPS-denied environments, multi-sensor fusion
071 algorithms have been developed for reliable odometry estimation (Khattak et al., 2019).

072 Research has also addressed depth visualization for navigation and obstacle avoidance in low-vision
073 scenarios (Lieby et al., 2011). Synthetic Vision Systems and full-windshield Head-Up Displays have
074 been explored to aid drivers and pilots in low visibility (Kramer et al., 2014; Charissis & Papanasta-
075 siou, 2010). Novel image enhancement methods for low-light conditions have been proposed (Atom
076 et al., 2020), as well as combinations of visual cues with standard wireless communication for road
077 safety (Boban et al., 2012). The importance of geometrical shapes and colors in Head-Up Displays
078 for driving perception has been emphasized (Zhan et al., 2023).

079 These advancements, however, face common challenges. These include increased computational
080 complexity due to sophisticated algorithms (Zhang et al., 2018; Atom et al., 2020; Tang et al.,
081 2022), durability and performance issues under variable or extreme environmental conditions (Khat-
082 tak et al., 2019; Boban et al., 2012), potential over-fitting problems due to limited datasets (Zhang
083 et al., 2018; Khattak et al., 2019), and the need for extensive real-world testing (Liu et al., 2014;
084 Boban et al., 2012; Tang et al., 2022). Some studies also lack clarity in explanations or comprehen-
085 sive validation (Kramer et al., 2014; Zhan et al., 2023).

086 In the realm of visual recognition and object detection, researchers have explored integrating human-
087 like processing mechanisms with computational models. Studies have delved into brain mechanisms
088 for object recognition, emphasizing hierarchical, feedforward processes (DiCarlo et al., 2012). Com-
089 parisons between human visual processing and deep neural networks (DNNs) have noted human
090 superiority in handling visual distortions and differences in attention mechanisms (Dodge & Karam,
091 2017; van Dyck et al., 2021). Attempts to direct DNNs’ visual attention using human eye-tracking
092 data have shown limited success in mimicking human attention patterns (van Dyck et al., 2022).

093 Innovative approaches include adversarial learning to enhance feature discrimination and match
094 feature priors (Yang et al., 2023a), biologically inspired models integrating top-down and bottom-
095 up processes for robust visual recognition in robotics (Malowany & Guterman, 2020), and models
096 mimicking the mammalian retina to enhance dehazing capabilities (Zhang et al., 2015). Some re-
097 searchers have proposed models using foveal-peripheral dynamics to reduce computational demands
098 while maintaining high-resolution perception in focused areas (Lukanov et al., 2021).

099 Recent studies have addressed specific challenges in low-visibility conditions such as fog, low light,
100 and sandstorms. The YOLOv5s FMG algorithm has been introduced for small target detection,
101 integrating various modules for better accuracy and localization (Zheng et al., 2023). Networks
102 improving image clarity in hazy and sandstorm conditions have been developed using novel MLP-
103 based modules for pixel reconstruction (Gao et al., 2023). The Prior Knowledge-Guided Adversarial
104 Learning (PKAL) approach leverages adversarial learning and feature priors for robust visual recog-
105 nition under adverse visibility (Yang et al., 2023b).

106 Enhancements to existing models, such as YOLOv8, have incorporated deformable convolutions
107 and attention mechanisms for better pedestrian and vehicle detection in poor visibility (Wu & Gao,
2023). Comprehensive reviews of image de-hazing techniques have highlighted limitations of non-

108 learning and meta-heuristic methods in real-time applications (V et al., 2023). The impact of low-
109 level vision techniques on high-level visual recognition tasks has been evaluated, suggesting a more
110 integrated approach for better outcomes in poor visibility conditions (Yang et al., 2020).

111 Novel techniques like spatiotemporal attention detection have been introduced to discern region-
112 level attention in video sequences (Zhai & Shah, 2006). The Parallel Detecting and Enhancing
113 Models (PDE) framework aims to simultaneously improve object detection and image enhancement
114 (Li et al., 2022). Research in visual saliency detection has explored integrating spatial position pri-
115 ors with background cues (Jian et al., 2021), while studies on early visual cues have examined their
116 role in detecting object boundaries in natural scenes (Mély et al., 2016).

117 Despite the advancements in the field, several limitations persist across existing studies. Many
118 approaches still struggle with joint optimization of object detection and image enhancement, the
119 detection of out-of-focus and low-contrast objects, and maintaining performance in dynamically
120 changing visibility conditions. This paper addresses these challenges by proposing a methodology
121 that combines human visual cues with computational models for object detection in low-visibility
122 conditions. By leveraging insights from human perception, such as attention mechanisms and con-
123 textual understanding, the proposed approach aims to enhance the robustness and accuracy of object
124 detection systems. This integration not only helps in effectively handling varying degrees of visi-
125 bility but also reduces computational complexity by focusing processing power on areas of interest,
126 similar to human selective attention.

127 Current techniques often struggle with the computational burden of processing high-resolution im-
128 ages in their entirety and may lack robustness in dynamically changing visibility conditions. Addi-
129 tionally, uniform dehazing techniques attempt to improve visibility across the entire image, which
130 can unnecessarily process visually clear regions. This leads to increased computational load and
131 potential degradation of image parts where high clarity is not essential. Our methodology addresses
132 these issues by focusing processing power on selected regions, reducing unnecessary computations.
133 Furthermore, by adapting the processing intensity based on real-time feedback, we enhance system
134 responsiveness and accuracy under diverse operational conditions.

135 The proposed approach stands out due to its unique integration of human visual cues into the object
136 detection process, particularly in low-visibility conditions. Unlike existing methods that may not
137 fully optimize computational resources or adapt to varying environmental conditions effectively, the
138 proposed architecture mimics the human eye’s capability to focus on relevant areas dynamically.
139 This method not only enhances detection accuracy but also improves computational efficiency by
140 prioritizing resource allocation, which is crucial for real-time applications. By addressing these key
141 aspects, this research aims to push the boundaries of object detection in low-visibility conditions,
142 offering a more robust, efficient, and adaptable solution compared to existing methods.

143 3 METHODOLOGY

144
145
146 The proposed methodology in Figure 1 focuses on developing a novel deep-learning framework in-
147 spired by the atmospheric scattering model and the human visual cortex to enhance object detection
148 in low-visibility conditions. The framework employs adaptive image enhancement techniques in-
149 tegrated with an object detection network to explore different integration strategies. The pipeline
150 initiates with a lightweight object detection model to identify regions of interest, which are subse-
151 quently leveraged for spatial attention in the dehazing process, followed by a more robust detection
152 model for refined and comprehensive object detection. This architecture will be evaluated across
153 various configurations using both synthetic and real-world foggy datasets, with performance mea-
154 sured using standard object detection metrics such as mean Average Precision (mAP) and image
155 quality metrics like Structural Similarity Index Measure (SSIM)A.1 and Peak Signal-to-Noise Ratio
156 (PSNR)A.2.

157 3.1 DATASETS

158 3.1.1 FOGGY CITYSCAPES

159
160
161 The Foggy Cityscapes (Sakaridis et al., 2018) dataset is created to address the problem of semantic
foggy scene understanding (SFSU). While there has been extensive research on image dehazing and

semantic scene understanding with clear-weather images, SFSU has received little attention. Due to the difficulty in collecting and annotating foggy images, synthetic fog is added to real images depicting clear-weather outdoor scenes. This synthetic fog generation leverages incomplete depth information to create realistic foggy conditions on images from the Cityscapes dataset, resulting in Foggy Cityscapes with 20,550 images. The training set consists of 2975 images, validated on 500 images, and the test set had 1525 images. The key features of the dataset include:

- **Synthetic Fog Generation:** Real clear-weather images are used, and synthetic fog is added using a complete pipeline that employs the transmission map.
- **Data Utilization:** The dataset can be utilized for supervised learning and semi-supervised learning. We generated a foggy dataset using the synthetic transmission map and then performed supervised learning on the synthetic foggy data.

3.1.2 RESIDE- β

The RESIDE- β Outdoor Training Set (OTS) is a comprehensive dataset designed to facilitate research in outdoor image dehazing. It addresses the challenges posed by haze in outdoor scenes, which significantly degrades image quality and affects subsequent tasks such as object detection and semantic segmentation. The dataset includes approximately 72,135 outdoor images with varying degrees of haze, enabling robust training of dehazing algorithms. For testing, we are using the RESIDE- β (REalistic Single Image DEhazing) dataset (Li et al., 2019). The subset of RESIDE- β , Real-Time Testing Set (RTTS) consists of 4,322 real-world hazy images with annotations for object detection. The training set consists of 3000 images, validated on 500 images, and the test set had 1500 images.

3.2 HUMAN VISUAL CUES

Selective Attention and Foveation: The human eye is not equally sensitive to all parts of the visual field. Central vision, or foveal vision, is highly detailed and used for tasks like reading and identifying objects. Peripheral vision is less detailed and more sensitive to motion. The system scans the entire image (peripheral vision) similar to our preliminary detection phase. This will identify areas for a more detailed analysis (foveal vision), mimicking the human approach of not processing every detail with equal clarity but focusing on areas of interest.

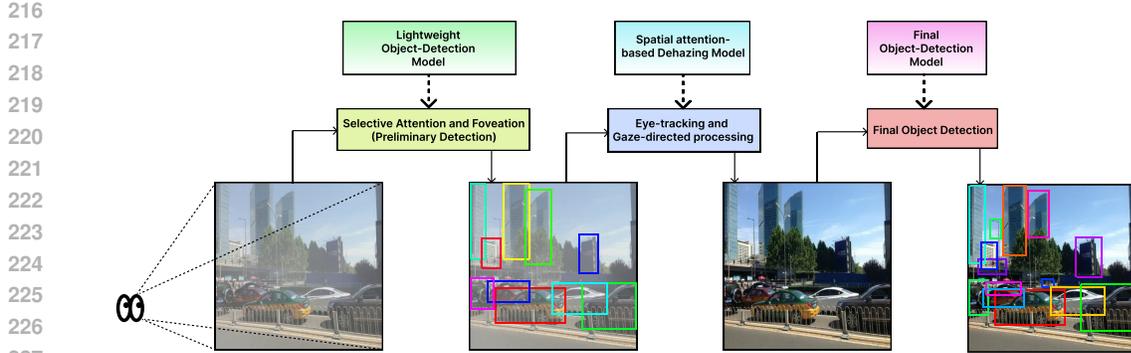
Adaptation to environmental conditions: Just as the human visual system adapts to different lighting conditions and levels of visibility (such as adjusting to a dark room after being in bright sunlight), the adaptive dehazing method adjusts the intensity and focus of its processing based on the detection feedback and environmental context, analogous to the way human vision adjusts to ensure optimal perception under varying conditions.

Eye Tracking and Gaze-directed processing: Eye-tracking monitors where a person is looking (the gaze) and what draws attention. In visual processing, this is analogous to directing computational resources toward areas of interest, much like the proposed method focuses on dehazing and detailed detection of regions where objects are likely to be present. By analogy, the system pays more attention to certain parts of the image, just as a person would fixate on specific areas within their field of view when searching for something.

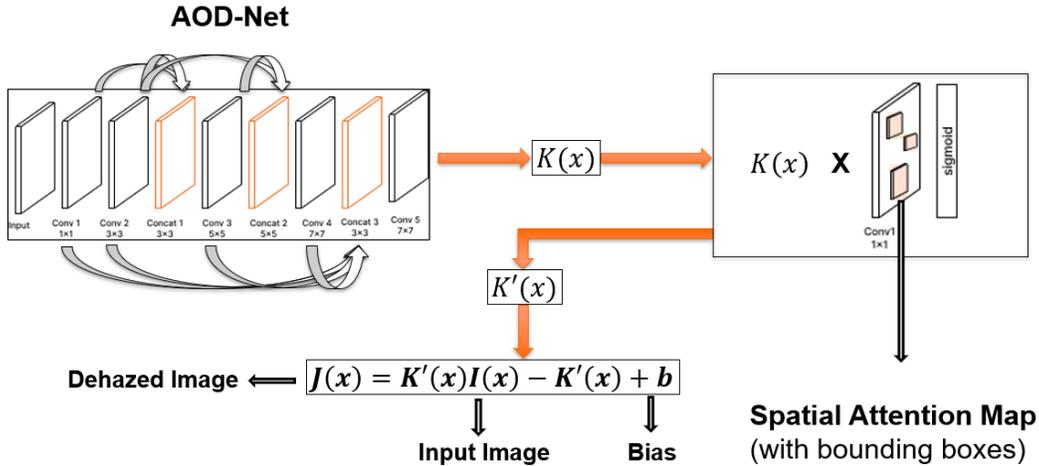
Integration of Bottom-up and Top-down processes: The human visual system uses both bottom-up signals (from sensory input) and top-down processes (based on knowledge, expectations, and current goals) to interpret scenes. The proposed model initially uses a bottom-up approach (object detection algorithms flagging potential areas of interest) followed by a top-down approach (focusing on dehazing efforts based on three flagged areas and previous learning), mirroring the complex interplay between sensory data and cognitive processes in human vision.

3.3 DEHAZING

Preliminary Detection: Implement a lightweight, fast object detection algorithm such as YOLOv5s or YOLOv8n to quickly scan the image for potential regions of interest or active regions and flag those patches with a high likelihood of containing objects. The smaller versions of standard YOLO



229 Figure 1: Overall architecture of Perceptual Piercing: (a)Preliminary detection using lightweight
230 object detection model (b) Gaze-directed dehazing using spatial attention on region of interests (c)
231 Final detection using a large and robust model



249 Figure 2: Architecture of AOD-NetX: It takes the transmission map output from the AOD-Net and
250 applies the spatial attention layer to focus on major areas of interest in the given input image.

251
252
253 models are less accurate than their full-sized counterparts but significantly faster, making them ideal
254 for preliminary detection.

255 **Region-based dehazing:** Apply dehazing algorithms specifically to those active regions identified
256 in the preliminary detection phase. Considering the depth or level of haze, the method should adapt
257 based on the characteristics of the detected regions

258 The proposed architecture of **AOD-NetX** in Figure 2 utilizes the transmission map created by the
259 standard AOD-Net (Li et al., 2017) and applies it within a spatial attention map module to produce an
260 attention-focused transmission map. This spatial attention map is derived from the bounding boxes or
261 Regions of Interest identified by the lightweight model (YOLOv5s/YOLOv8n) in our proposed
262 method. A sigmoid layer follows, mapping the output probabilities to a range between 0 and 1. We
263 opt not to use softmax in this context due to the independent significance of each bounding box.

264 3.4 OBJECT DETECTION MODELS

265
266
267 The YOLO models used in the detection pipeline include a variety of versions optimized for differ-
268 ent purposes. YOLOv5s is a lightweight variant designed for real-time detection with low compu-
269 tational requirements, while YOLOv8n (Nano) is tailored for high-speed applications on resource-
constrained devices like mobile phones. On the other hand, YOLOv5x, with its CSP backbone

Table 1: Performance of dehazing methods: AOD-Net and AOD-NetX

| Dataset | Dehazing Method | Evaluation Metrics | |
|----------------------|-----------------|--------------------|--------------|
| | | SSIM | PSNR |
| Foggy Cityscapes | AOD-Net | 0.994 | 26.74 |
| | AOD-NetX | 0.998 | 27.22 |
| RESIDE- β OTS | AOD-Net | 0.920 | 24.14 |
| | AOD-NetX | 0.945 | 25.80 |
| RESIDE- β RTTS | AOD-Net | 0.932 | 27.59 |
| | AOD-NetX | 0.656 | 27.62 |

and advanced data augmentation, provides enhanced performance for more complex scenes, and YOLOv8x (Extra Large) offers maximum accuracy for large-scale datasets. The detection process begins by using YOLOv5s or YOLOv8n on foggy images to generate initial annotations. These annotations, along with the original image, are then dehazed using AOD-NetX, and the dehazed image is subsequently passed through YOLOv5x or YOLOv8x for precise and refined detection results.

4 RESULTS

The dehazing modules are trained separately on the provided datasets, while the object detection models (various YOLO versions) remain as pre-trained on the MS-COCO dataset. This approach allows users to integrate the dehazing module with their own detection pipeline without requiring a complete re-training of the entire system. However, for improved results, the entire architecture could be fine-tuned on the target datasets, which would serve as a valuable direction for future ablation studies.

4.1 DEHAZING PERFORMANCE

The results in Table 1 show that AOD-NetX generally outperforms the standard AOD-Net in terms of SSIM and PSNR across most datasets. For Foggy Cityscapes and RESIDE- β OTS, AOD-NetX achieves higher SSIM and PSNR, indicating improved structural similarity and signal quality. However, for RESIDE- β RTTS, while AOD-NetX has a slightly better PSNR, AOD-Net achieves a significantly higher SSIM score, suggesting that AOD-Net may retain more structural details in this particular dataset. Overall, AOD-NetX is more effective in most scenarios, especially for complex foggy conditions.



(a) Foggy Cityscapes: Before Dehazing



(b) Foggy Cityscapes: After Dehazing (using AOD-NetX)

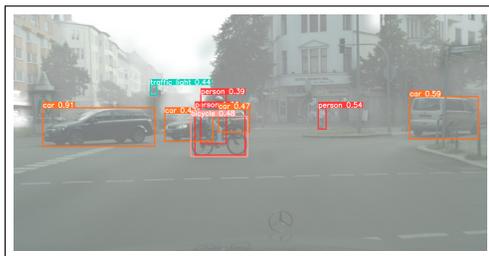
Figure 3: Dehazing performance on Foggy Cityscapes dataset

Table 2: **Train-** Foggy Cityscapes, **Test-** Foggy Cityscapes: Evaluation of various Perceptual Piercing variations based on mean Average Precision (mAP) under both clear and foggy conditions.

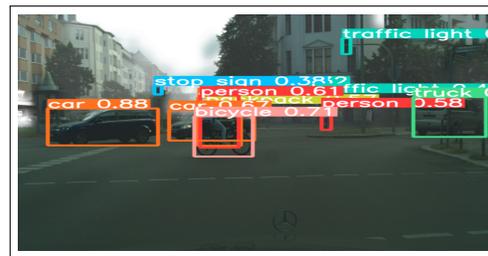
| Architecture Variants | Conditions | Evaluation Metrics (mAP) |
|--------------------------|------------|--------------------------|
| YOLOv5x | Clear | 0.5644 |
| | Foggy | 0.485 |
| AOD-Net+YOLOv5x | Clear | 0.6813 |
| | Foggy | 0.5822 |
| YOLOv5s+AOD-NetX+YOLOv5x | Clear | 0.4896 |
| | Foggy | 0.6152 |
| YOLOv8x | Clear | 0.5243 |
| | Foggy | 0.4948 |
| AOD-Net+YOLOv8x | Clear | 0.6099 |
| | Foggy | 0.5900 |
| YOLOv8n+AOD-NetX+YOLOv8x | Clear | 0.5150 |
| | Foggy | 0.6114 |

4.2 PERFORMANCE OF PERCEPTUAL PIERCING

The evaluation results of Perceptual Piercing variations in Table 2 trained and tested on the Foggy Cityscapes dataset indicate that integrating dehazing modules, such as AOD-Net and AOD-NetX, consistently improves object detection performance in both clear and foggy conditions. The ‘AOD-Net + YOLOv5x’ variant achieved the highest mAP under clear conditions (0.6813), while ‘YOLOv5s + AOD-NetX + YOLOv5x’ and ‘YOLOv8n + AOD-NetX + YOLOv8x’ demonstrated the best performance in foggy scenarios, with mAP scores of 0.6152 and 0.6114, respectively. In comparison, baseline YOLO models (YOLOv5x and YOLOv8x) showed lower detection accuracy, highlighting the significance of using enhanced dehazing techniques for better object detection in low-visibility environments.



(a) Foggy Cityscapes: Before Dehazing



(b) Foggy Cityscapes: After Dehazing (using AOD-NetX)

Figure 4: Dehazing performance on Foggy Cityscapes dataset

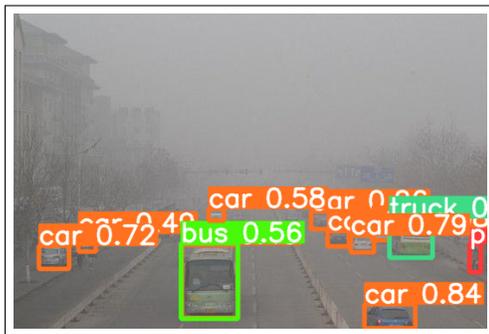
4.3 OUT-OF-DISTRIBUTION PERFORMANCE OF PERCEPTUAL PIERCING

The evaluation of various Perceptual Piercing variations in Table 3 trained on Foggy Cityscapes and tested on RESIDE- β OTS and RTTS datasets shows that the YOLOv8x architecture achieved the highest mAP scores under foggy conditions, with 0.7125 on OTS and 0.6978 on RTTS. Among the YOLOv5 variants, the baseline YOLOv5x model performed the best, with 0.6944 on OTS and

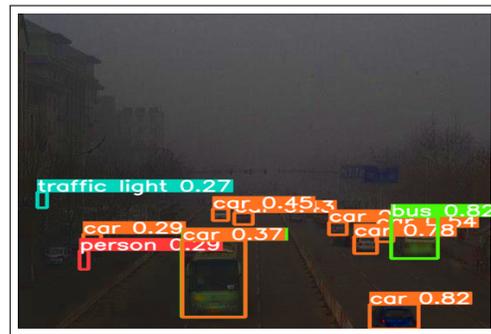
Table 3: **Train-** Foggy Cityscapes, **Test-** RESIDE- β OTS and RTTS: Evaluation of various Perceptual Piercing variations based on mean Average Precision (mAP) under foggy conditions.

| Architecture Variants | Configuration | Evaluation Metrics (mAP) |
|--------------------------|---------------|--------------------------|
| YOLOv5x | Test: OTS | 0.6944 |
| | Test: RTTS | 0.6655 |
| AOD-Net+YOLOv5x | Test: OTS | 0.6325 |
| | Test: RTTS | 0.6156 |
| YOLOv5s+AOD-NetX+YOLOv5x | Test: OTS | 0.5679 |
| | Test: RTTS | 0.5297 |
| YOLOv8x | Test: OTS | 0.7125 |
| | Test: RTTS | 0.6978 |
| AOD-Net+YOLOv8x | Test: OTS | 0.6458 |
| | Test: RTTS | 0.6125 |
| YOLOv8n+AOD-NetX+YOLOv8x | Test: OTS | 0.5779 |
| | Test: RTTS | 0.5312 |

0.6655 on RTTS. The addition of AOD-Net generally improved performance for YOLOv8 but had a diminishing effect on YOLOv5. Models incorporating AOD-NetX showed lower mAP values across both test datasets, indicating that its integration may need further optimization. Overall, the results suggest that YOLOv8x is more robust for foggy conditions compared to other variations.



(a) RESIDE- β : Before Dehazing



(b) RESIDE- β : After Dehazing (using AOD-NetX)

Figure 5: Dehazing performance on RESIDE- β dataset

5 DISCUSSION

Integrating a lightweight model with dehazing techniques forms a robust framework that significantly enhances overall system efficiency and effectiveness. This combined approach not only addresses the inherent limitations found in isolated systems but also synergizes their strengths to improve image clarity and object detection accuracy. By adopting a human-vision-inspired architecture, this methodology not only meets but exceeds the performance benchmarks set by state-of-the-art (SOTA) object detection models when tested against the same dataset distribution.

Furthermore, our directed dehazing strategy, which systematically targets specific image impairments, yields superior results with considerably fewer computations compared to traditional dehaz-

ing methods. This efficiency is pivotal, especially in real-time applications where computational resources and response times are critical factors. The success of this approach illustrates the potential of leveraging domain-specific enhancements to refine the capabilities of general object detection frameworks, suggesting a promising direction for future research and development in image processing technologies.

5.1 LIMITATIONS

The primary limitations of this paper are as follows: First, the proposed bio-inspired architecture does not incorporate image understanding from various low-visibility scenarios, which could have provided a more comprehensive validation of the methodology. Second, the scope of low-visibility images used is limited to foggy conditions, excluding other challenging environments such as rainy or hazy scenes. Extending the evaluation to rain or combined distribution datasets would enhance the robustness of the framework. Third, while the methodology aims for computational efficiency, the two-tiered detection process coupled with intensive region-specific dehazing may still require substantial computational resources, potentially limiting its applicability in real-time scenarios. Finally, in Out-of-Distribution (OOD) testing, the performance degrades compared to a more generalized model (e.g., YOLOv5x or YOLOv8x). It has been observed that even in clear images within the OOD dataset, the performance declines. This occurs because the dehazing model’s embedding space predominantly consists of foggy images, making it less effective when applied to clear scenarios in OOD datasets.

5.2 FUTURE WORK

To address the issue of generalizability in single-dataset training, two potential approaches are proposed. The first involves selectively applying the dehazing pipeline only when the scene is sufficiently hazy, determined using a simple haze index computed based on image contrast, brightness, and texture. The second approach is to train the dehazing model with embeddings from both foggy and clear images, thereby enabling it to generalize more effectively across diverse visibility conditions. To further enhance the robustness and applicability of our model, future research should focus on expanding testing with additional datasets that encompass a broader spectrum of low-visibility scenarios, including diverse environmental conditions such as rain, snow, and various levels of nighttime darkness. Such enhancements will enable the model to handle a wider range of adverse weather conditions, increasing its versatility and applicability in real-world situations. Moreover, incorporating training on more diverse datasets is crucial for improving generalization and optimizing performance in out-of-distribution testing. The availability of 4K datasets, which allow for the use of bounding box crops in dehazing, presents an opportunity to refine the model’s effectiveness further. Future efforts could also explore optimizing the model architecture and employing more advanced computational techniques to reduce resource demands, thereby enhancing feasibility for real-time applications in autonomous vehicles and other critical systems.

6 CONCLUSION

In conclusion, our research addresses the challenge of object detection under adverse conditions like fog, smoke, and haze, which commonly impair autonomous driving, aviation, and security. These environmental factors significantly degrade detection system performance, highlighting the need for precise, reliable methodologies. Our method uses a lightweight algorithm to identify regions of interest, followed by targeted dehazing to enhance visibility where needed most. The clarified images are processed through a robust detection model, boosting accuracy. This approach improves system efficiency and reliability for critical applications across various environments.

Our proposed AODNetX architecture outperforms state-of-the-art models, excelling in both standard and out-of-distribution datasets. This achievement aims to set new benchmarks in detection accuracy and efficiency. Moreover, our approach integrates atmospheric scattering model concepts and human visual cortex insights into machine learning frameworks. The expected outcome is an effective enhancement of object detection under challenging visibility, advancing safety and efficiency in technology-dependent sectors. This integration not only advances current detection systems but also deepens our understanding of visual processing in complex scenarios.

REFERENCES

- 486
487
488 Y. Atom, M. Ye, L. Ren, Y. Tai, and X. Liu. Color-wise attention network for low-light image en-
489 hancement. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Work-*
490 *shops (CVPRW)*, pp. 2130–2139, Seattle, WA, USA, 2020. doi: 10.1109/CVPRW50498.2020.
491 00261.
- 492 M. Boban, T. T. V. Vinhoza, O. K. Tonguz, and J. Barros. Seeing is believing—enhancing message
493 dissemination in vehicular networks through visual cues. *IEEE Communications Letters*, 16(2):
494 238–241, February 2012. doi: 10.1109/LCOMM.2011.122211.112093.
- 495 Vassilis Charissis and Stylianos Papanastasiou. Human–machine collaboration through vehicle
496 head up display interface. *Cognition, Technology Work*, 12:41–50, 2010. doi: 10.1007/
497 s10111-008-0117-0.
- 499 J. J. DiCarlo, D. Zoccolan, and N. C. Rust. How does the brain solve visual object recognition?
500 *Neuron*, 73(3):415–434, 2012. doi: 10.1016/j.neuron.2012.01.010.
- 502 S. Dodge and L. Karam. A study and comparison of human and deep learning recognition perfor-
503 mance under visual distortions. In *2017 26th International Conference on Computer Communi-*
504 *cation and Networks (ICCCN)*, 2017. doi: 10.1109/icccn.2017.8038465.
- 505 Y. Gao, W. Xu, and Y. Lu. Let you see in haze and sandstorm: Two-in-one low-visibility enhance-
506 ment network. *IEEE Transactions on Instrumentation and Measurement*, 72:1–12, 2023. doi:
507 10.1109/TIM.2023.3304668.
- 509 M. Jian, J. Wang, H. Yu, G. Wang, X. Meng, L. Yang, J. Dong, and Y. Yin. Visual saliency de-
510 tection by integrating spatial position prior of object with background cues. *Expert Systems With*
511 *Applications*, 168:114219, 2021. doi: 10.1016/j.eswa.2020.114219.
- 512 S. Khattak, C. Papachristos, and K. Alexis. Visual-thermal landmarks and inertial fusion for navi-
513 gation in degraded visual environments. In *2019 IEEE Aerospace Conference*, pp. 1–9, Big Sky,
514 MT, USA, 2019. doi: 10.1109/AERO.2019.8741787.
- 516 L. J. Kramer et al. Using vision system technologies to enable operational improvements for low
517 visibility approach and landing operations. In *2014 IEEE/AIAA 33rd Digital Avionics Systems*
518 *Conference (DASC)*, pp. 2B2–1–2B2–17, Colorado Springs, CO, USA, 2014. doi: 10.1109/
519 DASC.2014.6979422.
- 520 B. Li, X. Peng, Z. Wang, J. Xu, and D. Feng. Aod-net: All-in-one dehazing network. In *Proceedings*
521 *of the IEEE international conference on computer vision*, pp. 4770–4778, 2017.
- 523 B. Li et al. Benchmarking single-image dehazing and beyond. *IEEE Transactions on Image Pro-*
524 *cessing*, 28(1):492–505, 2019. doi: 10.1109/TIP.2018.2867951.
- 525 Zhiying Li, Shuyuan Lin, Zhongming Liang, Yongjia Lei, Zefan Wang, and Hao Chen. Pde: A real-
526 time object detection and enhancing model under low visibility conditions. *International Journal*
527 *of Advanced Computer Science and Applications(IJACSA)*, 13(12), 2022. doi: 10.14569/IJACSA.
528 2022.0131299.
- 530 P. Lieby et al. Substituting depth for intensity and real-time phosphene rendering: Visual navigation
531 under low vision conditions. In *2011 Annual International Conference of the IEEE Engineering*
532 *in Medicine and Biology Society*, pp. 8017–8020, Boston, MA, USA, 2011. doi: 10.1109/IEMBS.
533 2011.6091977.
- 534 C. Liu, Q. Zhao, Y. Zhang, and K. Tan. Runway extraction in low visibility conditions based on
535 sensor fusion method. *IEEE Sensors Journal*, 14(6):1980–1987, June 2014. doi: 10.1109/JSEN.
536 2014.2306911.
- 538 H. Lukanov, P. König, and G. Pipa. Biologically inspired deep learning model for efficient foveal-
539 peripheral vision. *Frontiers in Computational Neuroscience*, 15, 2021. doi: 10.3389/fncom.2021.
746204.

- 540 D. Malowany and H. Guterman. Biologically inspired visual system architecture for object recogni-
541 tion in autonomous systems. *Algorithms*, 13(7):167, 2020. doi: 10.3390/a13070167.
- 542
- 543 D. A. Mély, J. Kim, M. McGill, Y. Guo, and T. Serre. A systematic comparison between visual cues
544 for boundary detection. *Vision Research*, 120:93–107, 2016. doi: 10.1016/j.visres.2015.11.007.
- 545 C. Sakaridis, D. Dai, and L. Van Gool. Semantic foggy scene understanding with synthetic data.
546 *International Journal of Computer Vision*, 126:973–992, 2018. doi: 10.1007/s11263-018-1072-8.
- 547
- 548 Rong Tang, Qian Li, and Shaoen Tang. Comparison of visual features for image-based visi-
549 bility detection. *Journal of Atmospheric and Oceanic Technology*, 39, 2022. doi: 10.1175/
550 JTECH-D-21-0170.1.
- 551 S. Krishna B V, B. Rajalakshmi, U. Dhammini, M. K. Monika, C. Nethra, and K. Ashok. Image de-
552 hazing techniques for vision based applications - a survey. In *2023 International Conference for*
553 *Advancement in Technology (ICONAT)*, pp. 1–5, Goa, India, 2023. doi: 10.1109/ICONAT57137.
554 2023.10080156.
- 555 L. E. van Dyck, R. Kwitt, S. J. Denzler, and W. R. Gruber. Comparing object recognition in humans
556 and deep convolutional neural networks: An eye-tracking study. *Frontiers in Neuroscience*, 15,
557 2021. doi: 10.3389/fnins.2021.750639.
- 558
- 559 L. E. van Dyck, S. J. Denzler, and W. R. Gruber. Guiding visual attention in deep convolutional
560 neural networks based on human eye movements. *Frontiers in Neuroscience*, 16, 2022. doi:
561 10.3389/fnins.2022.975639.
- 562 X. Wu and Z. Gao. Based on the improved yolov8 pedestrian and vehicle detection under low-
563 visibility conditions. In *2023 2nd International Conference on Artificial Intelligence and In-*
564 *telligent Information Processing (AIIP)*, pp. 297–300, Hangzhou, China, 2023. doi: 10.1109/
565 AIIP61647.2023.00063.
- 566 J. Yang, J. Yang, L. Luo, Y. Wang, S. Wang, and J. Liu. Robust visual recognition in poor visibility
567 conditions: A prior knowledge-guided adversarial learning approach. *Electronics*, 12(17):3711,
568 2023a. doi: 10.3390/electronics12173711.
- 569
- 570 J. Yang, J. Yang, L. Luo, Y. Wang, S. Wang, and J. Liu. Robust visual recognition in poor visibil-
571 ity conditions: A prior knowledge-guided adversarial learning approach. *Electronics*, 12:3711,
572 2023b. doi: 10.3390/electronics12173711.
- 573 W. Yang et al. Advancing image understanding in poor visibility environments: A collec-
574 tive benchmark study. *IEEE Transactions on Image Processing*, 29:5737–5752, 2020. doi:
575 10.1109/TIP.2020.2981922.
- 576 Yun Zhai and Mubarak Shah. Visual attention detection in video sequences using spatiotemporal
577 cues. In *Proceedings of the 14th ACM international conference on Multimedia (MM '06)*, pp.
578 815–824, New York, NY, USA, 2006. Association for Computing Machinery. doi: 10.1145/
579 1180639.1180824.
- 580 Yu-Wei Zhan, Fan Liu, Xin Luo, Liqiang Nie, Xin-Shun Xu, and Mohan Kankanhalli. Generat-
581 ing human-centric visual cues for human-object interaction detection via large vision-language
582 models. *arXiv*, 2023. arXiv:2311.16475.
- 583
- 584 L. Zhang, Z. Zhai, L. Bai, Y. Li, W. Niu, and L. Yuan. Visual-inertial state estimation for the civil
585 aircraft landing in low visibility conditions. In *2018 International Conference on Networking*
586 *and Network Applications (NaNA)*, pp. 292–297, Xi’an, China, 2018. doi: 10.1109/NANA.2018.
587 8648726.
- 588 X.-S. Zhang, S.-B. Gao, C.-Y. Li, and Y.-J. Li. A retina inspired model for enhancing visibility
589 of hazy images. *Frontiers in Computational Neuroscience*, 9, 2015. doi: 10.3389/fncom.2015.
590 00151.
- 591
- 592 Y. Zheng, Y. Zhan, X. Huang, and G. Ji. Yolov5s fmg: An improved small target detection algorithm
593 based on yolov5 in low visibility. *IEEE Access*, 11:75782–75793, 2023. doi: 10.1109/ACCESS.
2023.3297218.

594 A EVALUATION METRICS

595 A.1 STRUCTURAL SIMILARITY INDEX MEASURE (SSIM)

596 The performance of dehazing methods is evaluated by the following equation of SSIM score between
597 two images:

$$598 \quad SSIM(x, y) = \frac{(2\mu_x\mu_y + c_1)(2\sigma_{xy} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x^2 + \sigma_y^2 + c_2)} \quad (1)$$

600 where:

- 601 • μ_x and μ_y are the average of x and y respectively.
- 602 • σ_x^2 and σ_y^2 are the variance of x and y respectively.
- 603 • σ_{xy} is the covariance of x and y .
- 604 • $c_1 = (k_1L)^2$ and $c_2 = (k_2L)^2$ are two variables to stabilize the division with weak de-
605 nominator; L is the dynamic range of the pixel-values (typically this is $2^{bits-per-pixel} - 1$),
606 $k_1 = 0.01$ and $k_2 = 0.03$ by default.

607 A.2 PEAK SIGNAL-TO-NOISE RATIO (PSNR)

608 Peak Signal-to-Noise Ratio (PSNR) is a widely used metric for evaluating the quality of recon-
609 structed images or videos compared to the original, reference data. It is expressed in decibels (dB)
610 and is calculated based on the mean squared error (MSE) between the original and the reconstructed
611 images. The formula for PSNR is given by:

$$612 \quad PSNR = 10 \cdot \log_{10} \left(\frac{MAX^2}{MSE} \right), \quad (2)$$

613 where MAX is the maximum possible pixel value of the image (for example, 255 for 8-bit images),
614 and MSE is the mean squared error, defined as:

$$615 \quad MSE = \frac{1}{mn} \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} (I(i, j) - K(i, j))^2, \quad (3)$$

616 where $I(i, j)$ represents the pixel value at position (i, j) in the original image, and $K(i, j)$ represents
617 the pixel value at the same position in the reconstructed image. Higher PSNR values generally
618 indicate better reconstruction quality, as they imply a lower MSE and thus less distortion. PSNR is
619 particularly useful for comparing the performance of different image processing algorithms in tasks
620 such as image compression, denoising, and super-resolution.

621 A.2.1 MEAN AVERAGE PRECISION (MAP)

622 For object detection performance, we are using mean Average Precision (mAP):

$$623 \quad AP = \frac{\sum_{k=1}^n (P(k) \times \text{rel}(k))}{\text{number of relevant documents}} \quad (4)$$

624 where:

- 625 • $P(k)$ is the precision at cutoff k in the list.
- 626 • $\text{rel}(k)$ is an indicator function equaling 1 if the item at rank k is a relevant document, 0
627 otherwise.
- 628 • n is the number of retrieved documents.

648 The mean Average Precision is then calculated as:

649

650

651

$$mAP = \frac{\sum_{q=1}^Q AP_q}{Q} \quad (5)$$

652

653

654

where AP_q is the Average Precision for the q^{th} query and Q is the total number of queries.

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701