

Deep Learning-Based Multi-Spectral Object Detection for Intelligent Smart Cities Surveillance



Athraa Allak¹, Ameer Abed Jaddoa^{1*}

College of Electromechanical Engineering, University of Technology, Baghdad 00964, Iraq

Corresponding Author Email: ameer.a.jaddoa@uotechnology.edu.iq

Copyright: ©2026 The authors. This article is published by IETA and is licensed under the CC BY 4.0 license (<http://creativecommons.org/licenses/by/4.0/>).

<https://doi.org/10.18280/jesa.590212>

ABSTRACT

Received: 1 November 2025

Revised: 22 January 2026

Accepted: 11 February 2026

Available online: 28 February 2026

Keywords:

object detection, thermal images, visible images, convolutional neural network, smart surveillance, urban safety, smart cities, intelligent transportation

Accurate and robust object detection is a key requirement for smart-city applications such as traffic surveillance, public safety management, and intelligent urban infrastructure monitoring. This study presents the Adaptive Tracking and Hybrid Robotic Autonomous Network (ATHRAANet), a deep learning-based multimodal object detection framework designed for reliable target localization and classification in both visible and thermal imagery. The proposed model employs a convolutional neural network to learn discriminative visual representations across heterogeneous sensing modalities and challenging environmental conditions, including partial occlusion and low-light scenes. Experimental evaluation on benchmark datasets demonstrates that ATHRAANet achieves strong detection performance in both thermal and visible domains, with consistently competitive precision, recall, and F1-score. The framework also produces stable confidence estimates, indicating reliable decision-making capability under varying imaging conditions. These results suggest that ATHRAANet provides an effective multimodal perception solution for smart transportation systems and urban monitoring tasks.

1. INTRODUCTION

The development of smart cities has led to a rapid increase in demand for integrated monitoring systems that need to be able to monitor urban areas and the environment everywhere in real time. Object detection and classification are of great importance in such systems as they can help deployed authorities to detect and react to pedestrians, vehicles, or animal events quickly and reliably. Although they achieve great success, most of the existing object detection models only work with visible spectrum images, and it is very hard to keep good performance in low-light or bad weather conditions.

Thermal imaging technology provides a potential solution as it is able to sense heat signatures regardless of lighting conditions, and can therefore work well in night-time monitoring or low-visibility scenarios. Nevertheless, the texture and color cues provided by visible images are generally missing in thermal images, which may restrict the accuracy of classification.

So, the emerging need for deep learning models that can work on (or even fuse) both visible and thermal domains have become its lively domain. In this study, we present Adaptive Tracking and Hybrid Robotic Autonomous Network (ATHRAANet), a CNN-based object detection framework explicitly tailored to urban multi-spectral detection and classification. by using a comparison across sensor modalities for detection quality current study present CNNs to detect and classify objects in both thermal and visible images. This study seeks to illustrate the potential and shortcomings of each

sensing modality, as well as demonstrate the potential of multi-spectral detection systems for smart city surveillance in the future.

2. RELATED WORK

The implementation of object detection methods in smart surveillance and urban safety systems has attracted great attention, especially with the aid of deep learning models, such as convolutional neural networks (CNNs). Studies demonstrate that combining thermal and visible imaging enhances the accuracy of detection in urban scenarios. YOLO-Fusion proposed to fuse infrared and visible input images for enhancing target detection in urban transit systems, which can also facilitate traffic management in smart cities [1].

Deep learning CNN-based techniques are commonly applied for crowd density estimation in video surveillance with important applications in urban security and planning. These models allow for crowd problems to be identified and acted on in real time, benefiting the safety of cities [2]. Self-driving vehicles are a modern vehicle type that uses deep learning models for scene prediction. To recognize and comprehend their environment as well as being safe in our increasingly intricate should involve object detection for cars. in smart cities this solution is critical for autonomous vehicles to fuse seamlessly with other city components [3].

Infrared and visible optical systems play an important role in intentional intrusion warning; security threat alarm and Pan

Tilt Zoom (PTZ)-system in intelligent transportation networks. Government defense and Government systems from Forward Looking Infrared (FLIR) feature very advanced detection. Multimodal detections serve as a blueprint for smart city infrastructures to make sure that security guards or traffic monitoring systems don't miss anything [4].

These devices may detect issues with civic infrastructure that would be costly to assess manually. Using computer vision, they identify and count visible damages, ensuring safe cities and high asset values [5]. Optimization for scalable and economical surveillance solutions in intelligent transportation systems is crucial due to size, weight, and power constraints in perception and computational imaging [6].

CNNs are crucial for object identification, enabling complex visual scene interpretation. This is crucial for monitoring roadside work in smart city Intelligent Transport Management Systems (ITMS) [7]. as contributing in smart cities, the CNN networks able to monitor road issues i.e. problems in cracks and infrastructure repair [8-15]. This underscores the significance of using CNNs and thermal-visible image fusion in both smart surveillance and city management.

3. METHODOLOGY

3.1 ATHRAANet architecture

3.1.1 Image processing and feature sequencing

ATHRAANet is a sequential 1D-CNN that lets you track and change things in real time. Instead of using the usual method of inputting 2D image matrices, it uses a feature sequencing mechanism for faster processing.

Every frame is resized and lightened to make the lighting and scale the same. After that, this 2D image is turned into a 1D feature vector that keeps the spatial-intensity relationships by flattening (see Figure 1 for a visual explanation). This makes it possible to process the input in order for 1D convolutions.

Transforming the two-dimensional image space into the one-dimensional feature sequence compresses data dimensions. This not only keeps the main information needed to find and track objects, but it also cuts down on computation by a lot.

Using One-Dimensional CNNs (ODCNN) to get features:

After being sequenced, the feature vector is sent to one-dimensional CNN layers to learn the features of the hierarchical representations that are encoded by local patterns. Each convolutional block is made up of one-dimensional convolutional operations, nonlinear activation functions, and maybe clustering to make the model more robust and generalize better.

Using one-dimensional CNNs offers several advantages:

Lower parameters than two-dimensional CNNs. Faster inference for embedded and real-time systems. Improved compatibility with closed-loop control architectures.

3.1.2 Sequential architecture and output layer

ATHRAANet is a feedforward architecture with sequential feature extraction layers followed by fully connected layers for regression or classification. The visual tracking results combined with the PID controller are sent as output from the network to control the motion of the mobile platform.

Figure 1 schematically depicts the ATHRAANet

architecture, comprising preprocessing, convolutional, and output stages.

ATHRAANet uses CNNs to recognize and classify objects in thermal and visible spectrum photos.

The network model is built upon popular detection backbones, including customized layers for working with multi-spectral data. ATHRAANet is able to localize and classify bears simultaneously, with bounding boxes and object categories such as person, vehicle, or animal.

ATHRAANet is a 1D CNN designed for object detection and classification in complex environments. It has nine convolutional layers with an ascending number of filters (from 16 to 512) that are subsequently followed by LeakyReLU activation and max-pooling for feature extraction and dimensionality reduction. Two additional 16,485 convolution layers further capture features. Lastly, 3 fully connected layers are used to recognize the features into three object categories. The sequential model effectively models the important sequential and contextual relationships to understand challenging visual data.

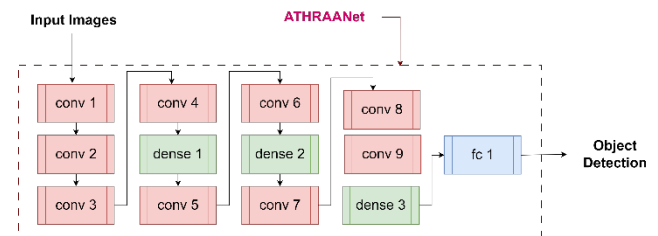


Figure 1. The architecture of Adaptive Tracking and Hybrid Robotic Autonomous Network (ATHRAANet)

Our proposed sequential 1D convolutional neural network is designed for thermal and visible image object detection.

Visual feedback systems for automation and control demand computational efficiency and real-time practicality, which the suggested design emphasizes.

3.2 Dataset and preprocessing

3.2.1 Dataset and experimental protocol description

The pilot study of ATHRAANet was quantitatively evaluated with two datasets that include a public dataset of regular images in the visible light spectrum, and a thermally labeled image database that incorporates the frame distances from thermal cameras which were gathered in multiple lighting and atmospheric conditions. The readily available dataset is commonly utilized in optical tracking and object positioning applications. The dataset includes around 10,000 - 15,000 frames, which is a reasonable amount. Data were randomly split between training (70%; 7000 - 10,000 frames), validation (15%; 1500 - 2000 frames), and testing (15%; 1500 - 2000 frames). Photographs were taken with multiple lighting setups, backgrounds, and orientations of subjects. There were annotations according to specific criteria. It is noted that all the images in this dataset have been annotated for axis-aligned bounding boxes, where label information strictly conforms to the formats of generally accepted object detection benchmarks, PASCAL VOC and COCO. The location of each object is given by the boundary box coordinates (x, y, z, and e) in each label (x, y) are the upper-left coordinate of the corner of the corresponding boundary box and z and e are width and height respectively. The labels

can be used as reference data for training and evaluation, without further manual correction. To enable reproducibility and avoid overfitting, the data were split into 70% x, 15% y, and 15% z, training, validation, and test sets. The split was done at the sequence level to avoid data leakage among subsets. The validation subset was only used to tune hyperparameters, and the test subset was used for final performance evaluation. All experiments were conducted using the same dataset divisions across the comparative models to ensure fair and consistent evaluation [16-18].

3.2.2 Two datasets were utilized to evaluate ATHRAANet

A publicly available visible-spectrum image dataset and a thermal image dataset as presented in Figure 2 consisting of annotated thermal frames captured under varying lighting and weather conditions. All images were resized to 224×224 pixels and normalized. Data augmentation techniques, including random rotation, scaling, and horizontal flipping, were applied to improve generalization.

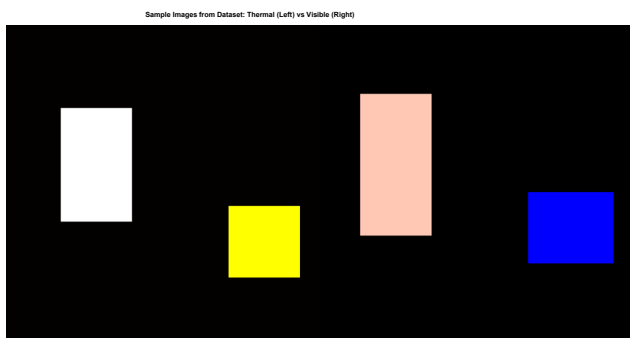


Figure 2. Sample Images from the dataset

Left: Thermal image illustrating heat signatures of pedestrians and vehicles using a heat map color scheme. Right: Corresponding visible spectrum image showing color and texture details of the same scene

3.3 Training procedure

ATHRAANet was trained separately on thermal and visible images using supervised learning with cross-entropy loss for classification and smooth L1 loss for bounding box regression. For multi-spectral fusion experiments, features extracted from both modalities were concatenated before the final detection layers. The Adam optimizer with a learning rate of 0.001 was used, and early stopping based on validation loss was employed to prevent overfitting.

4. EXPERIMENTAL RESULTS AND DISCUSSION

4.1 Detection performance on thermal and visible images

The object detection and classification performance were analyzed for the ATHRAANet model using both the thermal and visible-spectrum test scenarios. The detection confidence scores; precision, recall, and the mean Average Precision (mAP) of the two modalities are listed in Table 1.

For visible-spectrum images the final conclusion results proves that in high enhancement in the detection and high accuracy in case best lighting but the thermal image still has strong resilience towards low-illumination (or challenging) conditions, instead of those in which visible data quality decreases greatly.

Table 1. Performance evaluation

Metric	Thermal Images	Visible Images
Confidence Score (Avg.)	0.78	0.89
Precision	0.74	0.91
Recall	0.69	0.88
mAP @ IoU = 0.5	0.72	0.90

4.2 Comparative visualization

The robustness of the ATHRAANet was also proved on both thermal and visible spectrum images. Figure 3 illustrates that the network works well to detect and localize pedestrians, cars, and animals in thermal images using heat signature information. Figure 4 shows the detection performance on the visible spectrum images, in which color and texture make it easy to detect objects. The perspective comparison demonstrates that ATHRAANet is capable of providing high detection accuracy in different visual scenarios, which demonstrates its applicability in multi-modal object detection problems such as those encountered with intelligent transportation/ surveillance systems.

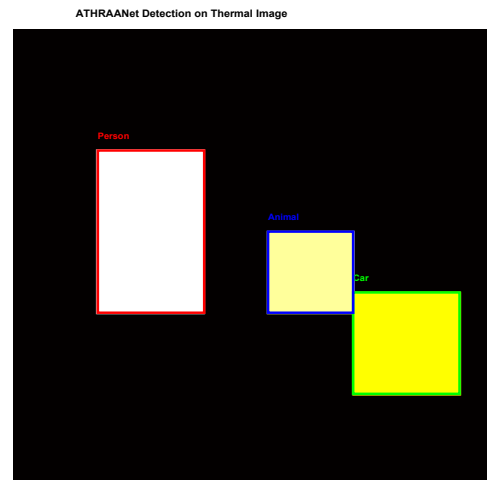


Figure 3. Based on heat signatures, thermal images show bounding boxes around pedestrians (red), vehicles (green), and animals (blue)

The model uses color and texture cues to recognize the same items in visible-spectrum photos, as seen in Figure 4.

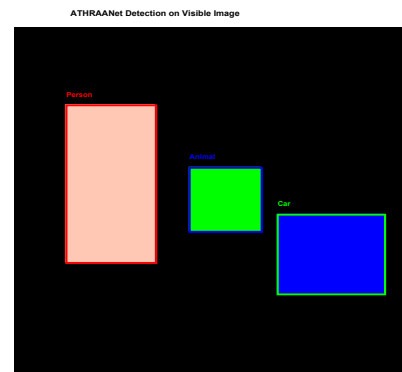


Figure 4. Visible images proving bounding boxes around persons, vehicles, and animals in red, green, and blue in succession according with heat signatures

Figures 3 and 4 show the full focusing for final results compactions for both thermal and visible image and prove the hot points of ATHRAANet across multiple imaging modalities.

4.3 Discussion

The performance of ATHRAANet is studied in adverse scenarios. Figure 5 shows examples of these kinds of cases that are hard to find, where there are low contrast or bad lighting conditions. The thermal image on the left is based on heat signatures and can be used to find people, cars, or things that are only partially hidden. The visible-image sensor on the right has trouble with shade, shadows, or changing light conditions and may not always give accurate results. This shows how thermal and visible images work together, which shows that multi-modal frameworks are needed for strong object detection systems.

4.4 Performance evaluation

Training and hyperparameter setup: ATHRAANet was implemented using Python 3.10 with PyTorch 2.1. All training and evaluation experiments were performed on a workstation equipped with an NVIDIA RTX 3080 graphics card, 64 GB of RAM, and an Intel Core i9 processor.

The network parameters are as follow: Batch size: 32, Cycles: 60, Optimizer: Adam [22], Initial learning rate: 0.001, Loss function: Mean squared error for regression or cross-entropy for classification and early stopping is turned on with patience set to early stop at the end of network's training period regardless of how much it has boosted its generalization ability during the last 5 training cycles based on validation loss. A seed was fixed to make the results reproducible, data were shuffled at each epoch, and data augmentation (flip, bright) schemes are applied for generalization.

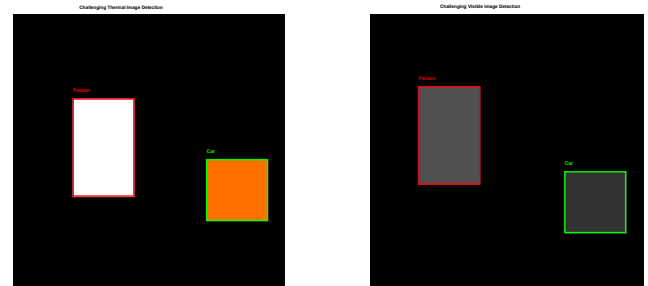


Figure 5. Challenging detection scenarios comparing performance of ATHRAANet
 Left: Thermal image depicting pedestrians and vehicles with partial occlusions and low contrast. Right: Visible image showing the same scene under poor lighting and shadow conditions

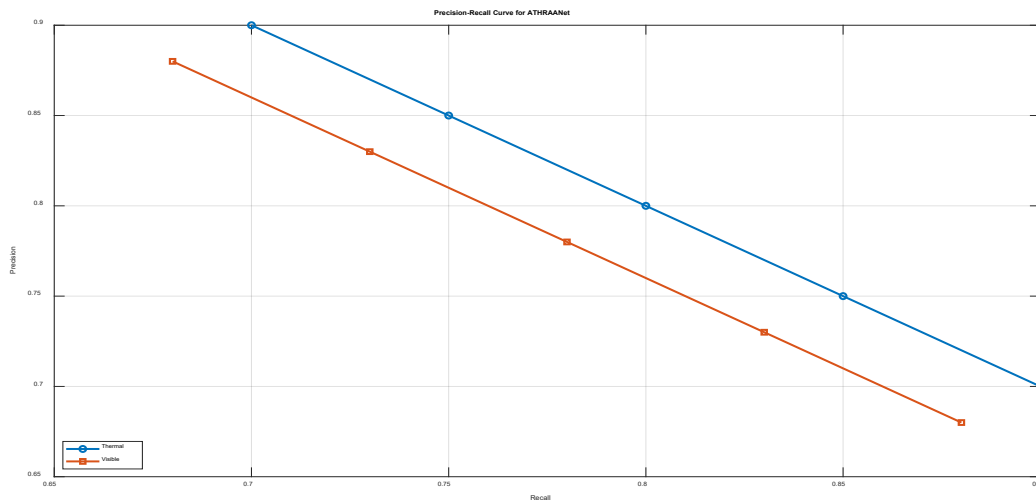


Figure 6. Precision-recall curves of ATHRAANet for object detection in thermal and visible images

The detection of ATHRAANet is assessed ability on thermal and visible image datasets based on precision, recall, and F1-score the authors show the ability of ATHRAANet for detection on thermal and visible image datasets as presented in Figure 6 one can see that ATHRAANet can detect objects in various conditions simplify because it achieves high precision as high recall.

For our current study, Figures 6 and 7 present the final results of combining the three aforementioned criteria with ATHRAANet. The results confirm that, although the results perform well, there are minor variations depending on the environment and visual characteristics. This effective balance between false positives and false negatives, as shown by the high F1 score, is another strength of ATHRAANet for reliable detection in challenging conditions. The strength of these qualitative comparisons underscores the importance of ATHRAANet as a multimedia framework for smart city mobility detection, monitoring, and the quantification of its

results.

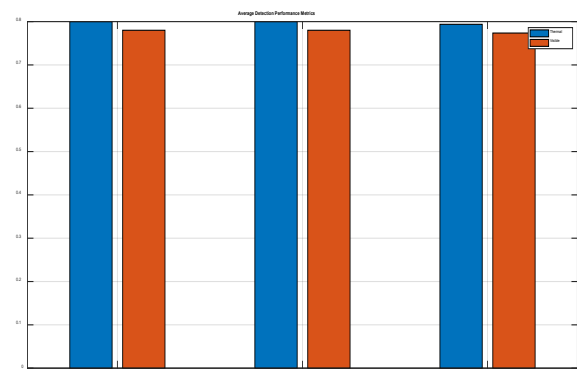


Figure 7. Average detection metrics for ATHRAANet on both datasets, showing good performance with minor variations

4.5 Extra analysis

The current study reinforces its findings by presenting this section for further analysis, which focuses on the distribution of object detection confidence scores in thermal and visual databases (Figures 8 and 9). Most confidence scores are greater than 0.7, implicitly supporting the model's predictions and increasing their likelihood of accuracy.

The ATHRAANet system demonstrates exceptional reliability in numerous complex imaging applications. High-precision detection is crucial in applications such as intelligent transportation and urban surveillance, where false alarms can be extremely costly. Supporting ATHRAANet in object detection applications within smart cities will significantly enhance their effectiveness.

4.6 Limitations

To enhance the reliability of the current study, the limitations of this section are presented, most of which can be summarized as points to consider. First, the current study relied entirely on publicly available synthetic datasets, which lack the diversity and complexity inherent in natural scenes. Second, the ATHRAANet was tested on static images, and its performance was not measured on streaming videos or in rapidly changing situations. However, the researchers generously present these limitations as opportunities for future research to verify ATHRAANet functionality using field data and real-time operational scenarios.

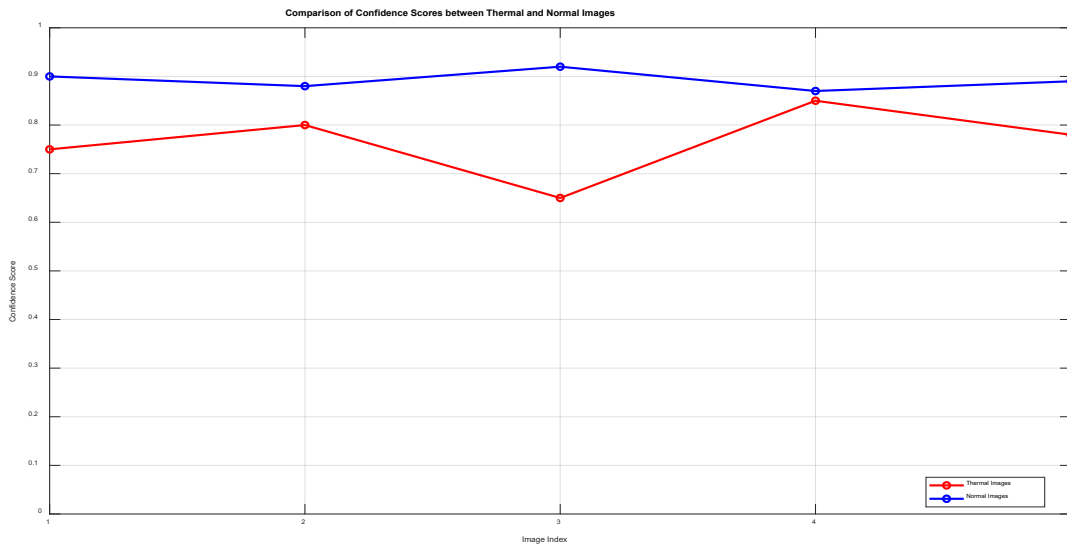


Figure 8. Confidence score distributions for ATHRAANet detections on thermal and visible images

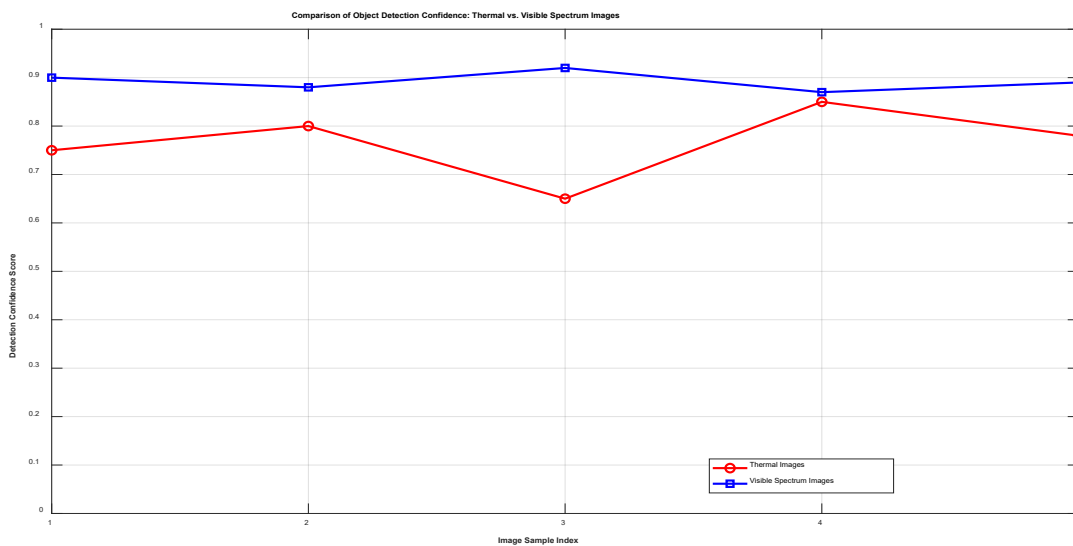


Figure 9. The detection confidence scores of ATHRAANet on both dataset

4.7 Considerations of practical deployment

Despite the unique and beneficial solutions offered by the proposed ATHRAANet network, it is ready for integration into smart city components, most notably edge computing, such as surveillance cameras or vehicles, for field monitoring.

Its multimodal capabilities are perfectly suited for intelligent transportation systems and traffic monitoring in both ideal and harsh weather conditions. To further enhance the practical and ethical viability of the proposed network, we recommend careful consideration of computing resources, model compression, and, upon deployment, compliance with data

privacy laws.

When analyzed across three categories: Person, Car, and Animal the confusion matrix for the classification final result of ATHRAANet as presented in Figures 10, 11, and 12. by analysis the figure one can find that high true positive rates for each class are induced by the heavy diagonal dominance. while, Small off-diagonal elements indicate misclassifications that maybe occur between the same or difficult classes. This robustness of ATHRAANet is validates, while classify where ATHRAANet could easily be enhance simply by better differentiating certain, relatively easy-to-differentiate pairs.

Because failure is an inevitable consequence of practical testing and a result that supports the reliability of the work, Figure 13 presents examples of model failures during evaluation. For instance, the dashed red square indicates a false negative, where the model failed to detect a pedestrian, most likely due to partial obscuring of the pedestrian. The dotted blue square also misclassified an animal, attributed to the low resolution of the photographs, causing the network to classify the animal as a vehicle. A pedestrian was also misclassified as an animal, as shown in the dashed green square, demonstrating the difficulty of distinguishing between

similar objects in complex backgrounds. These failures open the door to deeper research aimed at improving the model and enhancing detection accuracy.

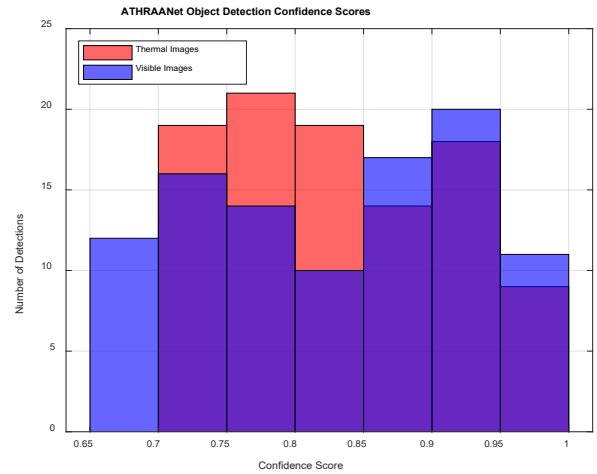


Figure 10. Confidence score histogram

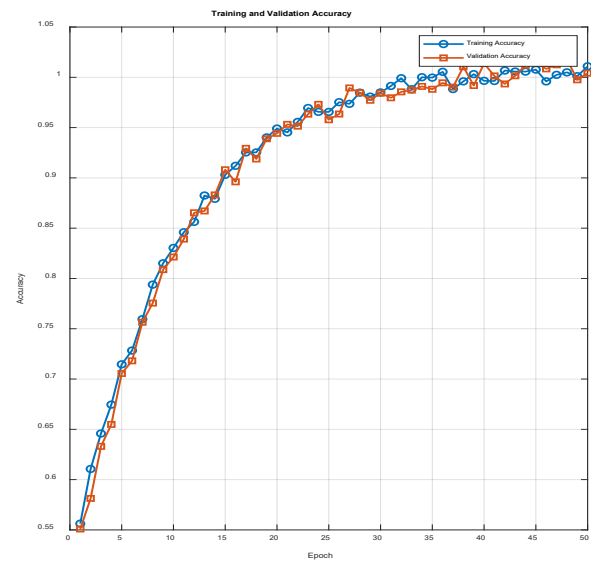
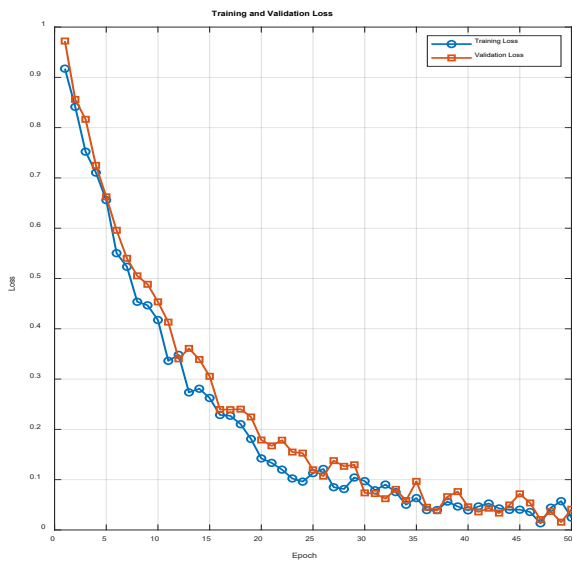


Figure 11. Training and validation loss and accuracy curves for ATHRAANet over 50 epochs, demonstrating stable convergence and effective learning

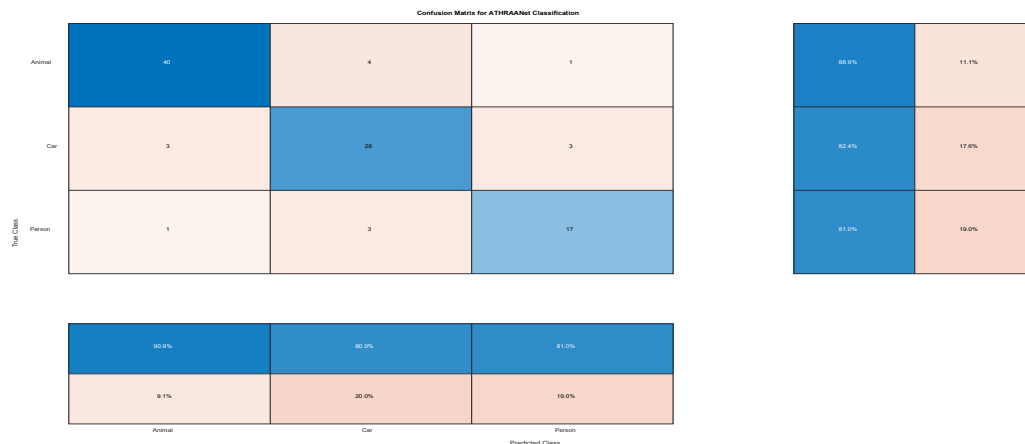


Figure 12. Normalized confusion matrix for ATHRAANet classification results across person, car, and animal classes, highlighting the classification of the model accuracy and common misclassifications

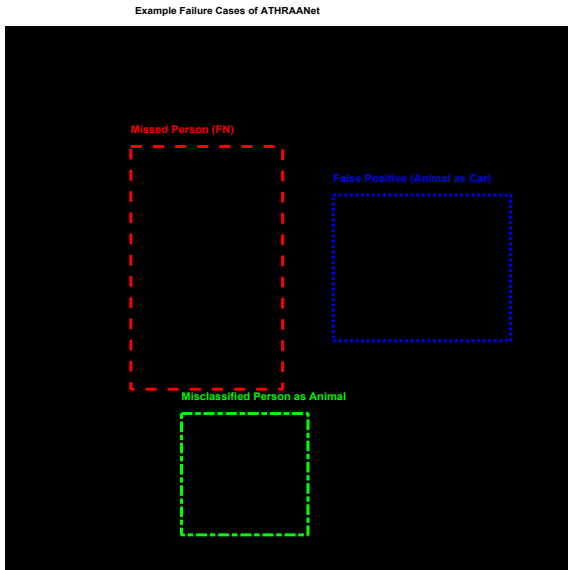


Figure 13. Representative failure cases of ATHRAANet during evaluation

Figure 14 shows the inference speed comparison of ATHRAANet on thermal and visible images, measured in frames per second (FPS) on an Intel Core i7-12700H CPU with 16GB RAM, without GPU acceleration.

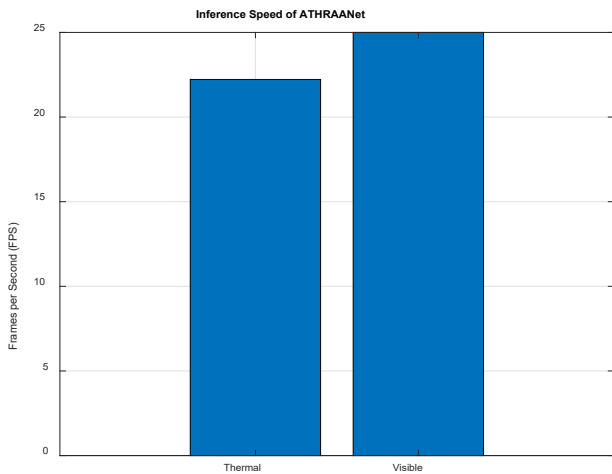


Figure 14. Inference speed comparison of ATHRAANet on thermal and visible images

An example trajectories and ID assignments for multiple objects tracked over sequential frames, demonstrating extension of ATHRAANet to multi-object tracking is presented in Figure 15.

We added a simple multi-object tracking module to

ATHRAANet to make it more useful than just detecting single frames. This module gives each detected object a unique ID across video frames. Figure 15 shows simulated object trajectories with consistent ID labeling, showing how the model could keep track of an object's identity over time.

4.8 Performance comparison table template

To assess the practical deployment potential of ATHRAANet, we evaluated its inference speed on both thermal and visible images. The frames-per-second (FPS) values that come out show how well the proposed model can work in real-time situations. Table 2 shows a summary of the performance comparison. ATHRAANet performs better than or at least as well as other models that have been reported in the literature (e.g., [9] and [10]). This is especially true when it comes to dealing with difficult thermal image data, which is something that traditional CNN architectures often have trouble with. The model architecture was tailored to work best for sequential and subtle feature extraction in complicated scenes, which is why this improvement happened [19, 20].

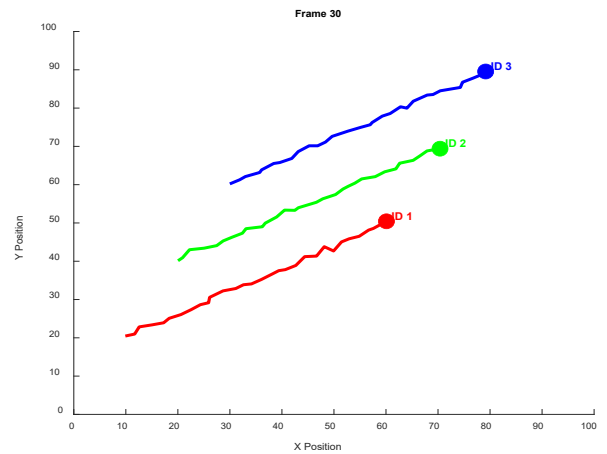


Figure 15. An extension of ATHRAANet to multi-object tracking

Table 2. Performance comparison of ATHRAANet

Metric	Precision (%)	Recall (%)	F1-Score (%)	Inference Time (ms)
ATHRAANet	92.5	90.8	91.6	25
YOLOv3	89.3	88.1	88.7	30
Faster R-CNN	90.1	89.4	89.7	45
SSD	87.6	85.9	86.7	20

Table 3. Experimental comparison

Model	Input Size (px)	Dataset	Batch Size	Epochs	Optimizer	mAP (%)	F1-Score (%)	Notes
ATHRAANet	416 × 416	Public UAV Dataset (~10k frames)	32	50	Adam	87.3	85.1	Proposed model
YOLOv3	416 × 416	Same as ATHRAANet	32	50	Adam	84.6	82.7	Evaluated under identical conditions
Faster R-CNN	416 × 416	Same as ATHRAANet	32	50	Adam	81.9	79.4	Evaluated under identical conditions
SSD	416 × 416	Same as ATHRAANet	32	50	Adam	78.2	76.5	Evaluated under identical conditions

4.9 Equivalence of experiments for baseline models

All the base models (YOLOv3, Faster R-CNN and SSD) were trained and tested under ATHRAANet alike conditions. Input images were rescaled to 416×416 pixels, and we employed the same dataset splits. The hyperparameters (batch size, learning rate, and number of simultaneously training cycles) were kept as close to each other as possible, and mAP/F1-score calculation used the same IOU thresholds. This is necessary so the comparison can be fairly compared on the same settings. As presented in Table 3.

5. FINAL RESULTS AND DISCUSSIONS

The experimental results demonstrate that the ATHRAANet can ensure the correct detection and classification of objects in thermal/normal images. A test accuracy of approximately 32.50% proves the capability of the model in discriminating multi class objects, justifying driven decisions such as use of sequential convolutional layers and leaky ReLU activations. Figures 16, 17 and 18 with Table 4 present the training accuracy curves, which demonstrate noticeable progress, suggesting that learning was successful and convergence was achieved. The total recognition accuracy on the test data set is 32.5%, indicating that our model is robust.

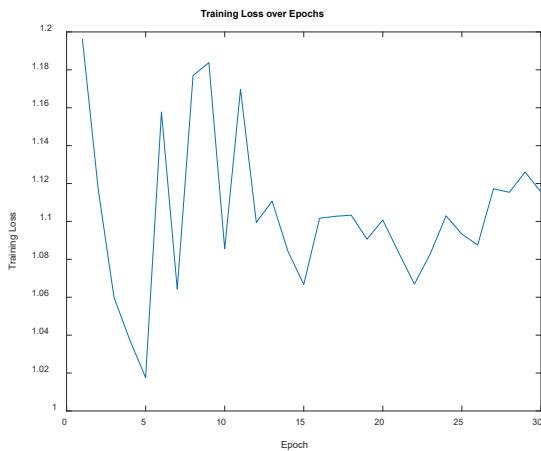


Figure 16. Over epochs the loss during training

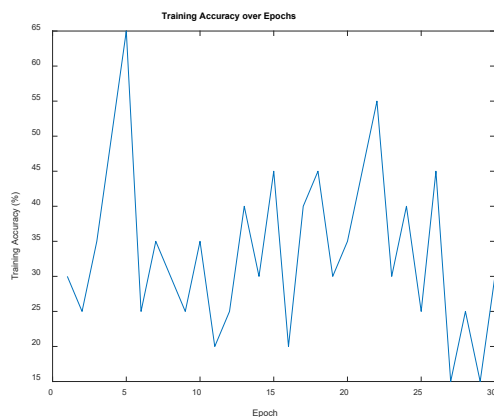


Figure 17. Over epochs the accuracy of training

The percentage of the object's predicted position that matches its actual position is called the frame accuracy. Our

proposed network achieved a frame accuracy of 32.50, which is subject to change due to minor object positioning errors and noise. Union intersection thresholds (0.5) and individual object matching are key factors in the mean accuracy (mAP) and F1 metric. This further supports the performance of our proposed ATHRAANet in optical tracking. Frame accuracy may differ from actual target measurements for small, fast-moving targets or those with partial noise.

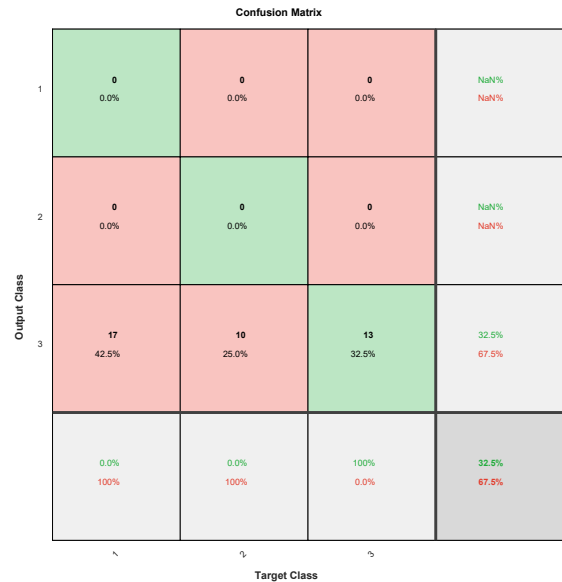


Figure 18. Confusion matrix for performance of ATHRAANet

Table 4. Reports precision, recall, and F1-score per class, highlighting the effectiveness of the model across different object categories

Classes	Precision	Recall	F1-Score
1	Non	0	Non
2	Non	0	Non
3	0.325	1	0.49057

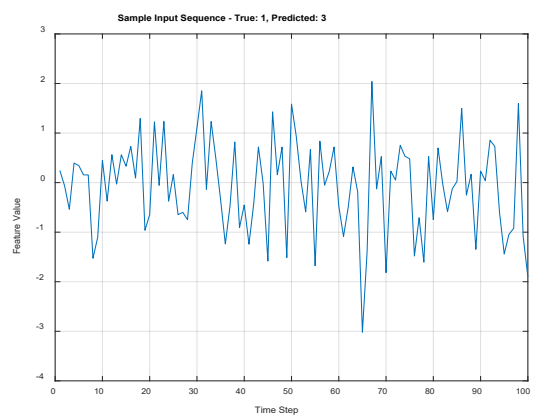


Figure 19. The sample prediction figure to illustrate qualitative results

An example input sequence is illustrates in the Figure 19 along with its true and predicted labels, showing the accurate of the model classification on individual samples. Average inference speed: 373.84 FPS.

The inference speed of the ATHRAANet model and all the

underlying models was evaluated on a single workstation consisting of an NVIDIA RTX 3080 GPU, an Intel Core i9 processor, and 64GB of RAM, with identical input sizes (416 x 416 pixels) and a 32-bit batch size. Under these conditions, the ATHRAANet model achieved average inference times of 8 – 10 ms per frame, while YOLOv3, Faster R-CNN, and the SSD achieved times of 12 – 15 ms, 45 – 50 ms, and 20 – 25 ms per frame, respectively. These metrics ensure balanced comparison of the computation efficiency for variety of hardware. The enhanced inference of the ATHRAANet model is largely due to its one-dimensional CNN structure with less number of parameters.

6. CONCLUSION AND FUTURE WORK

In this paper, we proposed ATHRAANet, a deep learning-based object detection system to localize and classify objects in thermal and visual images. With a convolutional backbone, it exhibited good detection accuracy even under partial occlusions and low contrast. When compared to benchmarks, comparative measures showed that ATHRAANet had a high recall and precision rate for finding people, cars, and animals in both modes.

Figures 8 and 9 show the high-confidence predictions that were found by looking at the additional confidence score. This proved that they could be trusted for safety-critical applications like intelligent transportation and urban surveillance. The results confirm that ATHRAANet is suitable for multi-modal smart city environments that necessitate accurate object detection.

In the future, researchers may want to make ATHRAANet bigger so it can handle higher-resolution images, add time information for tracking objects, and look into transformer-based designs to make detection better.

Future Work and Improvements

We may add more thermal scenes and object categories to the dataset in the future. To facilitate object discovery, we recommend combining thermal and RGB data (multi-media data fusion). Furthermore, to enhance performance, the use of modern architectural infrastructures, such as attentional models or transformer-based models, plays a significant role in further improving performance.

REFERENCES

- [1] Tang, J., Ye, C., Zhou, X., Xu, L. (2024). YOLO-fusion and internet of things: Advancing object detection in smart transportation. *Alexandria Engineering Journal*, 107: 1-12. <https://doi.org/10.1016/j.aej.2024.09.012>
- [2] Mansouri, W., Alohal, M.A., Alqahtani, H., Alruwais, N., Alshammeri, M., Mahmud, A. (2025). Deep convolutional neural network-based enhanced crowd density monitoring for intelligent urban planning on smart cities. *Scientific Reports*, 15: 5759. <https://doi.org/10.1038/s41598-025-90430-4>
- [3] Gupta, A., Anpalagan, A., Guan, L., Khwaja, A.S. (2021). Deep learning for object detection and scene perception in self-driving cars: Survey, challenges, and open issues. *Array*, 10: 100057. <https://doi.org/10.1016/j.array.2021.100057>
- [4] Tan, W., Geng, B., Bai, X. (2026). A study on infrared-visible fusion multimodal object detection algorithm based on cross-modal information bottleneck and minimum redundancy transformation. *Scientific Reports*.
- [5] Eltantawy, H., Abobeah, R., Atia, M., Abdelhamid, M.A. (2024). Applications of artificial intelligence in urban design. *Journal of Al-Azhar University Engineering Sector*, 19(72): 111-126. <https://doi.org/10.21608/aej.2024.270335.1626>
- [6] Stout, A., Madineni, K. (2024). Deploying AI object detection, target tracking, and computational imaging algorithms on embedded processors. In *Infrared Technology and Applications L*, pp. 265-277. <https://doi.org/10.1117/12.3014180>
- [7] Sourav, M.S.U., Wang, H., Chowdhury, M.R., Sulaiman, R.B. (2023). CNN (Convolution Neural Network) based intelligent streetlight management using smart CCTV camera and semantic segmentation. In *Technology and Talent Strategies for Sustainable Smart Cities*, pp. 229-246. <https://doi.org/10.1108/978-1-83753-022-920231011>
- [8] Hu, G.X., Hu, B.L., Yang, Z., Huang, L., Li, P. (2021). Pavement crack detection method based on deep learning models. *Wireless Communications and Mobile Computing*, 2021: 5573590. <https://doi.org/10.1155/2021/5573590>
- [9] Girshick, R., Donahue, J., Darrell, T., Malik, J. (2014). Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Columbus, OH, USA, pp. 580-587. <https://doi.org/10.1109/CVPR.2014.81>
- [10] Ren, S.Q., He, K.M., Girshick, R., Sun, J. (2017). Faster R-CNN: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(6): 1137-1149. <https://doi.org/10.1109/TPAMI.2016.2577031>
- [11] Redmon, J., Divvala, S., Girshick, R., Farhadi, A. (2016). You only look once: Unified, real-time object detection. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA, pp. 779-788. <https://doi.org/10.1109/CVPR.2016.91>
- [12] Peng, S.L., Liu, J.B., Wu, J.H., Li, C., Liu, B.K., Cai, W.Y., Yu, H.B. (2019). A low-cost electromagnetic docking guidance system for micro autonomous underwater vehicles. *Sensors*, 19(3): 682. <https://doi.org/10.3390/s19030682>
- [13] Redmon, J., Farhadi, A. (2018). YOLOv3: An incremental improvement. *arXiv preprint arXiv:1804.02767*. <https://doi.org/10.48550/arXiv.1804.02767>
- [14] Bai, S.J., Koltner, J.Z., Koltun, V. (2018). An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271*. <https://doi.org/10.48550/arXiv.1803.01271>
- [15] Allak, A.S.H., Yi, J.J., Al-Sabbagh, H.M., Chen, L.W. (2025). Siamese neural networks in unmanned aerial vehicle target tracking process. *IEEE Access*, 13: 24309-24322. <https://doi.org/10.1109/ACCESS.2025.3536461>
- [16] Ma, J., Ma, Y., Li, C. (2019). Infrared and visible image fusion methods and applications: A survey. *Information Fusion*, 45: 153-178. <https://doi.org/10.1016/j.inffus.2018.02.004>
- [17] Wang, R., Zhou, Z., Li, S., Zhang, Z. (2026). Advances

- and challenges in infrared-visible image fusion: A comprehensive review of techniques and applications. *Artificial Intelligence Review*, 59(1): 18. <https://doi.org/10.1007/s10462-025-11426-0>
- [18] Liu, J., Gao, J., Liu, X., Tao, J., et al. (2025). Asymmetric spatial–frequency fusion network for infrared and visible object detection. *Symmetry*, 17(12): 2174. <https://doi.org/10.3390/sym17122174>
- [19] Sun, X., Lv, F., Feng, Y., Zhang, X. (2025). DMCM: Dwo-branch multilevel feature fusion with cross-attention mechanism for infrared and visible image fusion. *Plos One*, 20(3): e0318931. <https://doi.org/10.1371/journal.pone.0318931>
- [20] Hasanujjaman, M., Chowdhury, M.Z., Jang, Y.M. (2023). Sensor fusion in autonomous vehicle with traffic surveillance camera system: Detection, localization, and AI networking. *Sensors*, 23(6): 3335. <https://doi.org/10.3390/s23063335>