

Pedestrian Detection in Low-Light Conditions: A Comprehensive Survey

Bahareh Ghari, Ali Tourani, Asadollah Shahbahrami, and Georgi Gaydadjiev

Abstract—Pedestrian detection remains a critical problem in various domains, such as computer vision, surveillance, and autonomous driving. In particular, accurate and instant detection of pedestrians in low-light conditions and reduced visibility is of utmost importance for autonomous vehicles to prevent accidents and save lives. This paper aims to comprehensively survey various pedestrian detection approaches, baselines, and datasets that specifically target low-light conditions. The survey discusses the challenges faced in detecting pedestrians at night and explores state-of-the-art methodologies proposed in recent years to address this issue. These methodologies encompass a diverse range, including deep learning-based, feature-based, and hybrid approaches, which have shown promising results in enhancing pedestrian detection performance under challenging lighting conditions. Furthermore, the paper highlights current research directions in the field and identifies potential solutions that merit further investigation by researchers. By thoroughly examining pedestrian detection techniques in low-light conditions, this survey seeks to contribute to the advancement of safer and more reliable autonomous driving systems and other applications related to pedestrian safety. Accordingly, most of the current approaches in the field use deep learning-based image fusion methodologies (*i.e.*, early, halfway, and late fusion) for accurate and reliable pedestrian detection. Moreover, the majority of the works in the field (approximately 48%) have been evaluated on the *KAIST* dataset, while the real-world video feeds recorded by authors have been used in less than six percent of the works.

Index Terms—Pedestrian Detection, Object Detection, Computer Vision, Autonomous Vehicles

I. INTRODUCTION

AUTOMATIC identification and localization of pedestrians in images or video frames captured by visual sensors have become increasingly vital in the computer vision domain. Pedestrian detection has use cases in various fields, such as autonomous vehicles [1], [2], surveillance systems [3]–[5], and robotics [6]–[8]. This task can be challenging to resolve in real-world scenarios, as there are different factors to consider for accurate performance. Accordingly, various illumination conditions, dissimilar pedestrian appearances and poses, occlusion, camouflage, and cluttered backgrounds can bring about issues for pedestrian detection systems [9]. Regarding the introduced challenges, low illumination is the leading problem

compared to others, as it has a natural or environment-related cause and cannot be prevented or naturally handled. It can be due to the time of the day, geographical location of where the scene is captured, weather conditions, *etc.* For instance, in some Scandinavian countries, especially during winter, the sunrise to sunset time can be less than ten hours and the pedestrian detection systems need to be adapted to the challenging scenarios.

Although some approaches have used Light Detection And Ranging (LiDAR) sensors, which are accurate remote sensing technologies that use laser light for distance measurement and obstacle avoidance, such sensors do not provide rich information from the surroundings [10]. There are many recently introduced pedestrian detection approaches such as [11]–[16] that cover the task in a wide range of scenarios. However, few recent works focus on detecting pedestrians at night and in low visibility conditions. In low-illumination scenarios, it is much more difficult for autonomous vehicles equipped only with vision sensors to detect moving objects on the road and prevent incidents. Fig 1 depicts how challenging the pedestrian detection task is compared to the same ordinary task during the daytime. The mentioned fact has resulted in an increase in demand for developing computer vision algorithms that can work under various illumination conditions.

Accordingly, this survey gives a deep review of 118 state-of-the-art low-light condition pedestrian detection algorithms. The research questions that the present work has aimed to answer are:

- **RQ1** Which datasets and baselines are mainly employed in low-light pedestrian detection tasks?
- **RQ2** What are the current deep learning-based algorithmic trends in low-illumination pedestrian detection?
- **RQ3** Regarding the state-of-the-art solutions, what are the currently existing and unresolved challenges in practical large-scale applications, such as fully autonomous vehicles?

To answer these questions, the paper in hand contributes to the body of knowledge in the field by providing contributions listed below:

- A survey of more than a hundred papers in the field of nighttime pedestrian detection,
- Review and classification of well-known baselines and datasets used for this purpose,
- Categorization of the state-of-the-art approaches in nighttime pedestrian detection regarding their architectural variations,

Bahareh Ghari and Asadollah Shahbahrami are with the Department of Computer Engineering, University of Guilan, Iran. baharehghari@msc.guilan.ac.ir, shahbahrami@guilan.ac.ir

Ali Tourani is with the Automation and Robotics Research Group, Interdisciplinary Centre for Security, Reliability, and Trust (SnT), University of Luxembourg, Luxembourg. ali.tourani@uni.lu.

Georgi Gaydadjiev is with the Computer Engineering Laboratory, Delft University of Technology, The Netherlands. g.n.gaydadjiev@tudelft.nl.

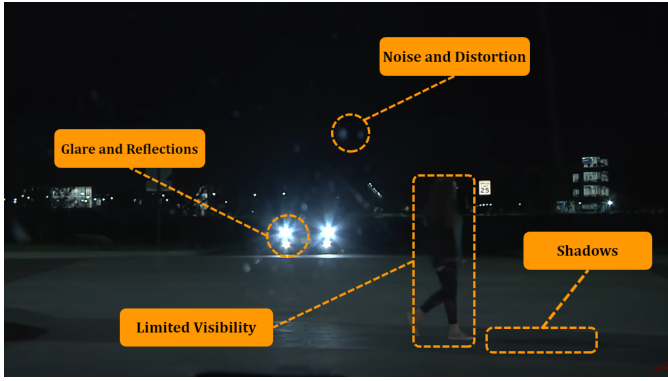


Fig. 1: Challenges of detecting pedestrians at night (image taken from 2019 Traffic Safety conference nighttime visibility report by Texas A&M Transportation Institute).

- Identification of the current trends and future methodologies in the field,

The rest of the paper is organized as follows: Section II reviews the currently available surveys in pedestrian detection in low-light conditions. Section III introduces the existing baselines and datasets employed by available approaches (RQ1). In Section IV, a set of recently introduced nighttime pedestrian detection approaches are introduced. Some discussions on the revealed trends and future insights in the field are presented in Section V (RQ2, RQ3). Finally, the paper concludes in Section VI.

II. RELATED SURVEYS

Given the significant importance of the pedestrian detection chore in cutting-edge domains such as autonomous vehicles, an extensive collection of surveys has been publicly available. These surveys delve into various aspects, including methodological approaches, target environmental contexts, evaluation procedures, and pre-defined presumptions. This section offers a concise study of these existing reviews and identifies the unexplored factors within them. This identification of gaps underscores the unique contribution that the manuscript in hand aims to provide in this field.

Chen *et al.* [17] analyzed various object detection methodologies along with robust feature extractors employed in the fields of vehicle and pedestrian detection. They also employed extensive experiments on the *KITTI* vision benchmark [18] as a well-known street dataset to assess the performance of the studied algorithms in terms of accuracy, inference time, memory consumption, model size, and the number of Floating-Point Operations per Second (FLOPS). It is important to note that their research primarily focuses on the algorithmic perspective within practical frameworks, addressing the algorithms' efficiency across diverse scenarios. Hou *et al.* [19] studied pixel-level image fusion strategies for vision-based pedestrian detection that works in all daytime/nighttime conditions and discussed efficient strategies of combining such methods with Convolutional Neural Network (CNN)-based fusion architectures. The primary aim of their research is to discuss a pixel-level fusion of strategies adopted from various

approaches that result in better performance for multi-spectral pedestrian detection tasks. Accordingly, their approach does not cover the future guidelines and possible strategies for such tasks. Authors in [20] studied deep learning-based methodologies employed in pedestrian detection tasks and provided informative discussions on how effective they are compared to other traditional algorithms. Although the mentioned survey also covers night-time pedestrian detection, the comparison among various methodologies was mainly established on different datasets with low-quality and multi-spectral instances. Other works such as [21], [22] surveyed approaches that targeted occlusion and scale variance challenges in pedestrian detection. They discussed solutions introduced in various papers that show acceptable performance in diverse conditions with occlusion, deformation, clutter, and scale difficulties.

Considering the introduced survey works, the present survey aims to provide specificities to set it apart from other available works, particularly regarding target scenarios and practical application use cases. In contrast with the previous works, the survey in hand is exclusively dedicated to nocturnal pedestrian detection, shedding light on the distinctive challenges tackled under low-light conditions introduced by state-of-the-art works. To the best of our knowledge, this work is the first study with a throughout focus on the detection of pedestrians in low-illumination (nighttime, in particular) conditions. In order to identify pedestrian detection approaches at nighttime that produce substantial results and feature novel architectures, the authors began by gathering and screening highly read and cited works from prominent venues over recent years. The sources included Google Scholar ¹, as well as well-established Computer Science bibliography databases, namely Scopus ² and DBLP ³. From the publications referenced in these sources, particular attention was given to those directly related to the ones targeted for challenging low-illumination conditions and underwent further checks to ensure their alignment with the domain. Following an in-depth exploration of the papers, they were systematically categorized based on their primary methodological solutions in addressing nighttime pedestrian detection challenges. Fig 2 depicts the distribution of papers collected, studied, and analyzed in the current survey.

III. BENCHMARKING DATASETS

Evaluation and development of pedestrian detection algorithms highly depend on providing proper data with annotated images/videos containing pedestrian instances. Such datasets should be well-annotated and cover diverse samples of pedestrian shots captured in real-world scenarios with various poses, occlusion levels, appearances, *etc.* to be considered appropriate for accurate training, testing, and validation stages. In this regard, this section collects various standard datasets for pedestrian detection at night that can be used for training and evaluation along with facilitating benchmark creation.

Ohio State University (OSU) Dataset ⁴ [31]: As one of the first pedestrian datasets, this thermal database provides a total

¹<https://scholar.google.com/>, accessed on 30 September 2023

²<https://www.dblp.org/>, accessed on 30 September 2023

³<https://www.scopus.com>, accessed on 30 September 2023

⁴<https://vcip1-okstate.org/pbvs/bench/Data/01/download.html>

TABLE I: Various datasets collected for pedestrian detection at night, sorted based on their publication year. Accordingly, dataset instances are collected using various sensors in different spectral ranges, *i.e.*, Near-Infrared (NIR), Middle-Infrared (MIR), and Far-Infrared (FIR).

Dataset	Metadata	Data						Sensor			
	Published	#Videos	#Frames	#Pedestrians	Resolution	Frame-rate*	Bit depth	RGB	NIR	MIR	FIR
OSU [31]	2005	10	~1.9k	984	360×240	30	8		✓	✓	✓
LITIV [39]	2012	9	~6.3k	-	320×240	30	8	✓	✓	✓	✓
CVC-09 [28]	2013	2	~11k	~14k	640×480	-	-				✓
LSI-FIR [34]	2013	13	~15.2k	~16.1k	164×129	-	14				✓
TIV [35]	2014	16	~63.7k	-	512×512	30	16		✓	✓	✓
KAIST [25]	2015	12	~95k	~103k	640×480	20	8	✓	✓	✓	✓
NTPD [38]	2015	-	~22k	-	64×128	-	-		✓	✓	✓
CVC-14 [29]	2016	4	~8.5k	~9.3k	640×512	10	-	✓			✓
KMU [33]	2016	23	~12.9k	-	640×480	30	24				✓
UTokyo [32]	2017	-	~7.5k	~2k	640×480	1	-	✓	✓	✓	✓
CAMEL [36]	2018	26	~43k	~80k	336×256	30	24	✓	✓	✓	
NightOwls [23]	2018	40	~279k	~42k	1024×640	15	-	✓			
SCUT [26]	2018	21	~211k	~477k	720×576	25	8				✓
YU FIR [40]	2018	-	~2.8k	~9.3k	640×480	30	14				✓
FLIR [41]	2020	-	~10.2k	~28.1k	640×512	24	16				✓
ZUT [30]	2020	-	~110k	~80k	640×480	30	16		✓	✓	✓
LLVIP [24]	2021	26	~33.6k	-	1080×720	1	24	✓	✓	✓	✓
C3I [37]	2022	6	~39k	-	640×480	30	8		✓	✓	✓

*presented in frames per second (fps)

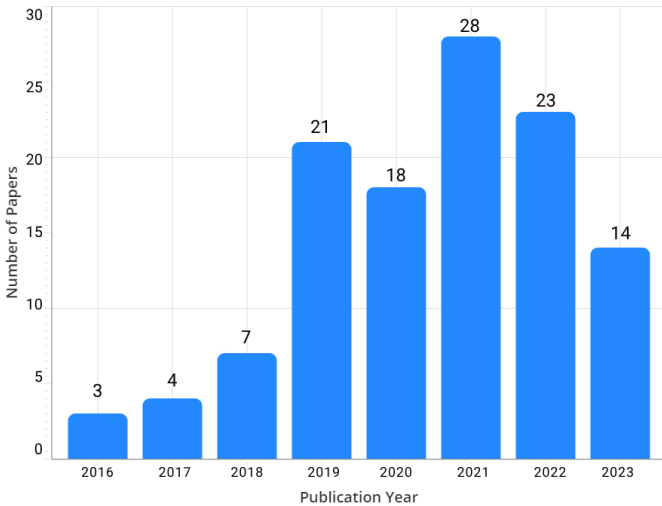


Fig. 2: Distributions of the papers surveyed in the current research work that only focus on pedestrian detection at low-light scenarios from 2016-2023 (total: 118 papers).

number of 1.9k thermal frames with a resolution of 360×240 , which were captured on campus and street. The *OSU* comprises three different classes of objects, including persons, cars, and poles. A total of 984 people were annotated in this dataset.

Laboratoire d’Interprétation et de Traitement d’Images et Vidéo (LITIV) ⁵ [39]: The dataset has nine video sequences, each containing people in an indoor hall with various zoom settings. The main challenges in these video sequences are the strong occlusions of objects and cluttered backgrounds.

CVC-09 Dataset ⁶ [28]: As another well-known dataset, CVC-09 is acquired during the day and night with 11k frames. The

dataset contains training and testing sets, where the day and night sequences contain 5,990 and 5,081 frames, respectively.

Laboratorio de Sistemas Inteligentes Far-Infrared (LSI-FIR) Dataset ⁷ [34]: This dataset is composed of classification and detection portions and contains grayscale images collected in different temperatures with varying illumination. The classification part has 16,152 positive samples (*i.e.*, pedestrian) and 65,440 negative samples (*i.e.*, background), while the detection part includes 15,224 images, categorized into 6,159 train and 9,065 test instances.

Thermal Infrared Video (TIV) ⁸ [35]: The dataset contains video sequences with 63,782 annotated frames for visual processing tasks, such as detection, counting, group motion estimation, and single-view and multiple-view tracking. Three out of sixteen sequences are mainly used for pedestrian detection, while other classes like car, runner, bicycle, and motorcycle are marked in this dataset.

Korea Advanced Institute of Science and Technology (KAIST) ⁹ [25]: It is one of the first multi-spectral pedestrian datasets with 95k aligned color-thermal image pairs and 103k dense annotation of samples. Data was captured from various traffic scenarios in the daytime and nighttime for autonomous driving applications. Annotations were manually added, resulting in three primary categories (person, people, and cyclist), and three occlusion levels (no-occlusion, partial-occlusion, and heavy occlusion).

Night-Time Pedestrian Dataset (NTPD) Dataset [38]: It contains a set of pedestrian images recorded by an active night vision system. The dataset contains 1,998 positive and 8,730 negative in the training set and 2,370 positive and 9,000 negative samples in the testing set.

⁷<https://www.kaggle.com/datasets/muhammedalkran/lsi-far-infrared-pedestrian-dataset/code>

⁸<http://csr.bu.edu/BU-TIV/>

⁹<https://github.com/SoonminHwang/rgbt-ped-detection>

⁵<https://www.polymtl.ca/litiv/en/codes-and-datasets>

⁶<http://adas.cvc.uab.es/elektre/enigma-portfolio/item-1/>



Fig. 3: Instances of some datasets introduced for nighttime pedestrian detection. It should be noted that the collected image or video sequences were captured using various sensors.

CVC-14 Dataset¹⁰ [29]: An extended version of the CVC-09 dataset, titled CVC-14, was introduced later to facilitate the challenges of automated driving. It contains video sequences of grayscale visible and thermal pairs corresponding to daytime and nighttime, where the daytime and the nighttime shares are 4,401 and 4,117 instances, respectively.

Keimyung University (KMU) Dataset¹¹ [33]: As a dataset captured using a FIR camera mounted on a vehicle driving in the summer nights for pedestrian detection, it contains three types of videos regarding the driving speed (20-30 km/h). It also covers pedestrians with different activities and poses, such as walking, running, and crossing the road. KMU has 4,474 positive and 3,405 negative frames in the training set and 5,045 frames in the testing set.

UTokyo¹² [32]: This multi-spectral dataset contains RGB, NIR, MIR, and FIR images collected in a university for object detection in automated driving, including person, car, and bike.

It contains 7,512 images, where 3,740 was taken during the daytime and the rest at nighttime.

CAMEL¹³ [36]: The dataset provides visible-infrared video sequences for multiple object detection and tracking, where 43k visible-infrared image pairs are annotated with four different object classes, including person, bike, vehicle, and motorcycle. CAMEL covers various real-world scenes, occluded targets, and different illumination conditions.

NightOwls¹⁴ [23]: This dataset targets the research on pedestrian detection at night and contains videos recorded in seven cities across Germany, the Netherlands, and the United Kingdom. It contains 279k frames with 42k pedestrians that have been manually labeled. Three primary labels (*i.e.*, far, medium, and near) have been assigned to the pedestrians to categorize them based on the distance they had from the vehicle during data acquisition. Additionally, frame brightness levels (low, medium, and high) and pedestrian pose (frontal and sideways) are other classification metrics employed in NightOwls.

¹⁰<http://adas.cvc.uab.es/elektro/enigma-portfolio/cvc-14-visible-fir-day-night-pedestrian-sequence-dataset/>

¹¹<https://cvpr.kmu.ac.kr/KMU-SPC.html>

¹²http://www.mi.t.u-tokyo.ac.jp/projects/mil_multispectral/

¹³<https://camel.ece.gatech.edu/>

¹⁴<https://www.nightowls-dataset.org/>

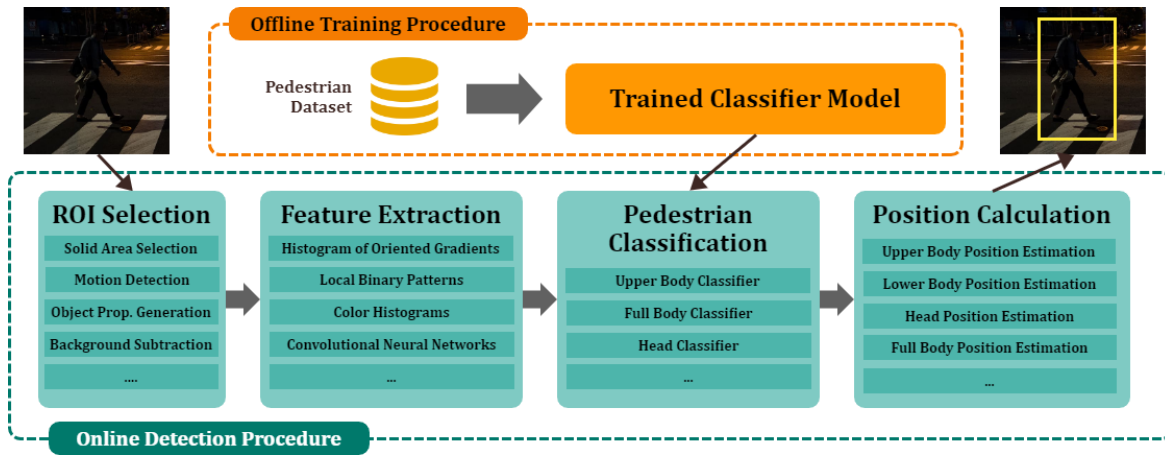


Fig. 4: The overall diagram of night-time pedestrian detection methodologies.

South China University of Technology (SCUT) Dataset¹⁵ [26]: A large-scale nighttime pedestrian dataset proposed by Xu *et al.* to motivate more attempts toward the task of on-road FIR pedestrian detection. The dataset contains approximately 11 hours-long image sequences with 211k annotated frames and a total of 477k bounding boxes for 7k unique pedestrians. SCUT groups pedestrians into three subsets, including near-scale (*i.e.*, ~ 80 pixels), medium-scale (*i.e.*, ~ 30 to ~ 80 pixels), and far-scale (*i.e.*, less than 30 pixels) subset based on the range of imaging distances.

YU FIR [40]: This seasonal temperature-based pedestrian detection dataset is captured on campus and urban traffic roads. The temperature was calibrated from -40°C to 150°C and used as the thermal infrared data for pedestrian detection. YU FIR contains a total of 2,802 frames with 1,803 and 575 positive images in the training set and test set, respectively.

Forward Looking InfraRed (FLIR) Dataset¹⁶ [41]: This multi-spectral dataset was collected for Advanced Driver Assistance Systems (ADAS) during daytime and nighttime. It contains visible-thermal image pairs, some of which are not aligned, and the rest contain 5k multi-spectral pairs for training and testing. This version contains three frequent object categories, including persons, bicycles, and cars.

Zachodniopomorski Uniwersytet Technologiczny (ZUT) Dataset¹⁷ [30]: It is a thermal dataset recorded in four European countries during diverse weather conditions, including sunny, foggy, heavy rain, light rain, and cloudy. The dataset contains 110k frames with 80k pedestrian annotations and provides synchronized Controller Area Network (CAN bus) data, including brake pedal status, driving speed, and outside temperature for ADAS.

Low Light Visible Image Person (LLVIP) Dataset¹⁸ [24]: The dataset is recorded by a binocular camera containing visible light and Infrared (IR) sensors. Targeting low-illumination surveillance tasks, the dataset contains 15k pairs of visible-infrared images. The annotations of IR and visible-light images

are the same due to the similar resolution and Field-of-View (FoV) of the cameras.

C3I Thermal Automotive Dataset¹⁹ [37]: The dataset was acquired in various environmental (*i.e.*, roadside, industrial town, alley, and downtown) and weather (*i.e.*, cloudy, foggy, windy, and sunny weather) conditions during daytime, evening, and nighttime. It comprises video sets with 39,770 frames, of which 17,740 frames are recorded in daytime, 12,640 in evening time, and 9,390 frames at nighttime. The frames are annotated in six object classes: person, car, bike, bicycle, bus, and pole.

To provide a comprehensive introduction to the datasets at hand for nighttime pedestrian detection, Fig. 3 depicts some of their instances. Additionally, Table I provides an in-depth description, facilitating a deeper understanding of their attributes and characteristics. It should be noted that these datasets have been curated in a way that encompasses a wide range of scenarios, including various lighting conditions, diverse pedestrian poses, occlusions, and complicated backgrounds. Such diversities ensure that provided data can serve as valuable resources for evaluating algorithms under real-world conditions.

IV. STATE OF THE ART

When considering the pedestrian detection methodologies for nighttime and low-illumination conditions, one of the primary architectures that come to mind for designing such frameworks is to include an *offline training* procedure that utilizes a dedicated pedestrian images dataset to train a classification model. In this regard, the model learns hidden patterns and characteristics specific to pedestrians in darker environments, enabling it to distinguish them from other objects or background elements. Although the mentioned architecture can be very beneficial, it should be noted that learning-based methodologies are not always used for this task. Fig. 4 shows these typical stages for pedestrian detection at night and how they are connected to each other. We can see the typical stages that form the foundation and serve as critical

¹⁵https://github.com/SCUT-CV/SCUT_FIR_Pedestrian_Dataset

¹⁶<https://www.flir.com/oem/adass/adass-dataset-form/>

¹⁷<https://iee-dataport.org/open-access/zut-fir-adass>

¹⁸<https://github.com/bupt-ai-cz/LLVIP/>

¹⁹<https://iee-dataport.org/documents/c3i-thermal-automotive-dataset/>

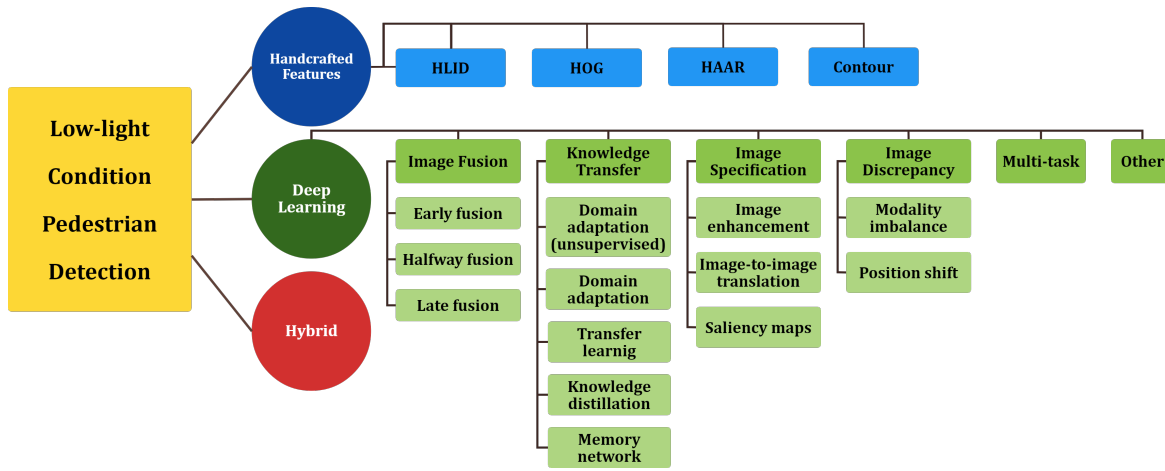


Fig. 5: The primary classification of different nighttime pedestrian detection methodologies considered in this survey.

components in identifying pedestrians, irrespective of the specific methodology employed. Whether the framework employs handcrafted features, machine learning models, or a mixture of both, it typically encompasses standard stages, including Region of Interest (ROI) selection (*i.e.*, identifying potential pedestrian regions within an image), visual feature extraction (*i.e.*, capturing relevant information from the selected ROIs), pedestrian classification (*i.e.*, using the features to classify the detected regions as pedestrians or non-pedestrians), and position calculation (*i.e.*, determining the precise location of the detected pedestrians).

In this survey, nighttime pedestrian detection approaches based on the underlying techniques and methodologies employed have been categorized into three distinct groups, including handcrafted features, deep learning, and hybrid methods. To understand the diverse strategies, along with their set of advantages and limitations, this section provides an in-depth study of the state-of-the-art works classifiable in the mentioned three categories.

A. Handcrafted Features Approaches

Handcrafted features contain manual design and selection of particular visual features from the input image/frame. Extracting information using such features has the advantage of simplicity, transparency, explainability, and the ability to provide consistent results across similar scenarios. They generally require lower computational costs and can still work well when there is no access to annotated data. However, they can have adaptability problems regarding their domain expertise and may reach a *performance ceiling*, making them difficult to be improved over a certain point. Considering these trade-offs, handcrafted features can still be found in many pedestrian detection frameworks under challenging illumination scenarios.

Regarding the intrinsic characteristics of handcrafted features, the majority of approaches in this category require employing thermal data. As one of the first works in this category, Davis *et al.* [31] used a combination of generalized person template derived from Contour Saliency Map (CSM) and background subtraction to identify pedestrians' locations

in thermal frames. Then, an *AdaBoost* classifier could validate the candidate regions, which adaptively adjusts the filters from the gradient information of training instances. While the template-based method brings about a quick screening procedure, it is considered a challenging methodology to detect groups of people in the scene. Similarly, Nowosielski *et al.* [42] presented a HAAR and Adaboost-based night-vision framework to identify humans in thermal images. The proposed algorithm processed all frames independently and without the aggregation mechanism, which increases the false positive rate due to incorrect recognition of the region as a person. An approach titled Thermal Infrared Radiometric Cumulative Channel Feature (TIR-ACF) introduced in [40] employs a thermal normalization methodology to factor in the maximum human body temperature for pedestrian detection. However, the experimental environment of this normalization strategy only includes a specific temperature range of small distant targets. As a more complicated methodology, Jeong *et al.* [33] presented an approach based on a Cascade Random Forest (CaRF) classifier, low-dimensional Haar-like features, and Oriented Center-Symmetric Local Binary Patterns (OCS-LBPs) for detecting sudden pedestrian crossing in thermal images. As the thermal temperature of the road is similar to or slightly higher than the pedestrians during summer nights, the concentration of this approach is on pedestrian samples in the summer season which leads to high prediction accuracy. Kim *et al.* [43] designed a pedestrian detector using a multi-level cascade learning algorithm and Histogram of Oriented Gradients (HOG) features. They used a smartphone-based thermal camera to capture human images of indoor environments to validate their work. Additionally, the 2D thermal image is mapped into a 3D space through an inverse perspective transformation method [44] to estimate the distance of the pedestrian detected from the camera.

Infrared images are another source of valuable information for pedestrian detection tasks at night. Zhou *et al.* [45] designed a pedestrian extraction algorithm for IR images. They build a global model using the weighted Histogram of Local Intensity Difference (HLID) and texture weighted HOG algorithms to locate potential pedestrian regions. Then,

using a head template based on the HAAR-like features and incorporating it into a local model for pedestrian head search, the global and head templates are combined to identify pedestrians. As another approach, Khalifa *et al.* [46] introduced a foreground detection framework that models the background's global motion between consecutive frames by applying the block-matching algorithms to the ROI to compensate for the camera motion. They use a Support Vector Machine (SVM) classifier to differentiate between the image's foreground and background. The evaluation results on the *CVC-14* show that the proposed algorithm can capture the dynamic aspect between frames in a video stream. Shahzad *et al.* [47] suggested a new procedure for pedestrian detection, tracking, and head detection in IR systems using template matching, Kalman filter, and HAAR cascade classifiers, respectively. The authors confirmed that the template matching method performs better than the contour-based method for pedestrian detection, and pedestrian tracking using the Kalman filter has the highest error rate. Likewise, and based on visual saliency in IR images, Cai *et al.* [48] proposed a model to focus on ROI generation along with a HLID feature and an SVM classifier to make a final detection. Considering that the visual saliency-based method includes small processing regions for candidate verification, the proposed algorithm demonstrates a fast execution time.

To conclude, the approaches with handcrafted features can provide acceptable results in many cases. However, they generally suffer from their incapability to handle complex scenarios due to their low discriminative nature and seem to act less flexibly while adapting to new scenarios.

B. Deep learning approaches

Solutions based on deep learning leverage the potential of neural networks to learn and extract features from raw image data automatically. In this regard, adaptability, versatility, generalization w.r.t. diverse scenarios, and high-performance results are among the expected outcomes of employing Deep Neural Networks (DNNs). They also have the capability to automatically learn features and reduce the requirement for manual feature engineering, along with integrating feature extraction and detection steps to have an end-to-end learning procedure. However, approaches in this category typically require large amounts of labeled data for training and powerful hardware due to their computationally intensive nature. Additionally, they lack straightforward explainability, leading to challenging interpretations of their decision-making process and, thus, fine-tuning to improve their performance.

Many recent works for low-light pedestrian detection employ DNNs as an inevitable part of their algorithms. These methodologies have been divided into categories below in this survey regarding their use case:

1) **Image Fusion Methodologies:** Image fusion refers to extracting and fusing the most significant characteristics of raw images captured by multiple sensors to generate a single image with complementary information, a compelling description of the scene, *etc.* [49]. Considering the *fusion* stage, CNN fusion architectures can be divided into three primary strategies,

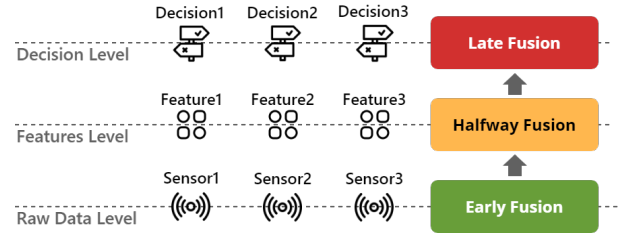


Fig. 6: Various image fusion strategies used in night-time pedestrian detection approaches.

namely, *early fusion*, *halfway fusion*, and *late fusion*. Fig 9 depicts a brief overview of different fusion architectures used in the research works.

Early fusion-based methods: in the context of night-time pedestrian detection, it indicates integrating visual and thermal feature maps right after the first *convolutional* layer of a CNN. However, fusing IR and RGB images to generate a four-channel input before feeding the network is another approach that can execute low-level feature fusion at an early stage. As a recent work introduced in [50], a You Only Look Once (YOLO) v3-based [51] multi-spectral pedestrian detector is introduced that can use RGB, thermal, and multi-spectral images. The mentioned approach merges the three channels of RGB and the single channel of thermal images to prepare a 4-channels input. The authors also evaluated a YOLO-4L [52] version to improve the detection accuracy of small-scale pedestrians on the scene. Evaluation results on various datasets demonstrated that the methodology outperforms other image types under all lighting conditions.

Halfway fusion-based methods: as a widely explored fusion strategy in recent years, the fusion operation of input modalities happens at the middle stages of a network, after the fourth convolutional layer. As one of the halfway fusion works, Yang *et al.* [53] designed a Cascaded Information Enhancement Module (CIEM) and a Cross-modal Attention Feature Fusion Module (CAFFM) to enrich the pedestrian information and suppress the interference caused by background noise in the color and thermal modalities. While CIEM uses a *spatial attention mechanism* to weigh the features combined by the cascaded feature fusion block, CAFFM employs *complementary features* to construct global features. In [54], Channel-wise Attention Module (CAM) and Spatial-wise Attention Module (SAM) were integrated into a multi-layer fusion CNN, aiming to re-weigh cross-spectral features at channel-dimension and pixel-level, respectively. Although the SAM methodology results in reduced detection speed, the performance of the approach is substantially improved. Zhang *et al.* [27] suggested a new halfway fusion strategy that applies cyclical fusion and refinement operations to achieve the consistency and complementary balance of multi-spectral features by controlling the number of loops. Based on the fact that the fused features are more discriminative than the mono-spectral ones, their main idea is to consecutively refine the spectral features with the fused ones and increase the overall feature quality. Hence, according to the analysis on the KAIST and FLIR dataset, the authors suggested that the

number of loops should be tuned for any dataset. A cross-modal framework based on YOLO v5 detector introduced in [55] for multi-spectral pedestrian detection. In their study, the information complementarity of RGB-thermal streams was acquired by a Cross-modality Feature Complementary Module (CFCM) to reduce the target loss. They use an Attention-based Feature Enhancement Fusion Module (AFEFM) to fuse different modalities' essential features and suppress the background noise while strengthening the semantic information. In another approach, and based on YOLO v5 lightweight network, Fu *et al.* [56] proposed an adaptive spatial and pixel-level feature fusion module, called *ASPPF Net*, to obtain fusion weights of spatial positions and pixel dimensions in two feature maps. The fusion weights are employed to recalibrate the original feature maps of visible and IR images to acquire multi-scale fusion feature layers. The spatial and pixel attention mechanisms enable the *ASPPF Net* to focus on learning useful information and suppress redundant information to achieve a fast prediction speed of 35 frames per second (fps) and lower Miss Rate (MR) on the night subset of the KAIST dataset. A Multi-Layer Fusion network based on Faster R-CNN (MLF-FRCNN) was proposed by [57], which employs Feature Pyramid Network (FPN) and Region Proposal Network (RPN) as two parallel feature extractors to deal with pedestrian samples with different scales. As a two-stage multi-spectral pedestrian detector, the MLF-FRCNN achieves a running time of 0.14 seconds per frame and the highest Average Precision (AP) in detecting various pedestrian scales. In [58], four variants of fusion models have been designed at different stages, titled low-level (*i.e.*, early fusion), middle-level (halfway fusion), high-level (late fusion), and confidence-level (score fusion). The first three approaches implement convolutional feature fusions, while the last corresponds to the combination of confidence scores from RGB and thermal CNN branches at the decision stage. The study reveals that the halfway fusion model achieves the lowest overall MR.

In halfway fusion spatial attention-based mechanisms, the importance of each location in the feature map is calculated to highlight the areas with valuable information. Accordingly, Cao *et al.* [59] used Channel Switching and Spatial Attention (CSSA) in a lightweight fusion module to effectively fuse multi-modal inputs while ensuring low computational cost. During channel switching, the channel of each modality with insufficient features is replaced by the corresponding channel from another modality. Likewise, a bi-directional fusion strategy called *BAA Net* is introduced in [60] to ensemble the RGB-thermal features for multi-spectral pedestrian detectors. The strategy distills the high-quality features of two modalities and re-calibrates the representations gradually. It contains intra- and inter-modality attention modules to improve spectra-specific features and adaptive selection of information from the most reliable modalities, respectively. In another similar work, Zhang *et al.* [61] introduced a two-stream CNN, titled Guided Attentive Feature Fusion (GAFF), to dynamically re-weight and integrate multi-spectral pedestrian features under the guidance of the intra- and inter-modality attention mechanisms. The intra-modality attention module aims to enhance the visible or thermal features in pedestrian areas, while

the inter-modality attention module selects the most reliable modality according to the feature quality, which requires costly annotation information. The authors' solution to this issue is to assign labels based on the prediction of pedestrian masks from the intra-modality attention module and then select the most relevant modality where the prediction mask is closer to the ground truth. Qingyun *et al.* [62] proposed a Cross-modality Fusion Transformer (CFT) module and embedded it to the YOLO v5 framework. The CFT learns long-range dependencies and focuses on global contextual information. In particular, by leveraging the self-attention mechanism, the network can simultaneously carry out intra-modal and inter-modal fusion and capture the latent interactions between visible and Thermal-Infrared (TIR) spectrums.

Some works discuss the most common feature fusion strategies in CNNs: *concatenation* (*i.e.*, stacking two feature maps at the exact spatial locations), *summation* (*i.e.*, calculating the sum of two feature maps at the exact spatial locations), *maximum* (*i.e.*, obtaining the maximum response of two feature maps at the exact spatial locations), and *mean* (*i.e.*, calculating the mean value of two feature maps at the exact spatial location) [63]. Accordingly, Pei *et al.* [63] discussed the influence of these strategies in various CNN fusion architectures, including merged Visual-Optical (VIS) and IR images based on *RetinaNet* detector [64]. The results proved that the summation fusion strategy performs better than other methodologies. Ding *et al.* [65] employed a Network-In-Network (NIN) in Region-based Fully Convolutional Network (R-FCN) framework [66] to merge the image information of two sub-networks to deal with large-scale and small-scale pedestrian instances. After the concatenation of Conv-VIS and Conv-IR, the small- and large-scale pedestrian candidates generated by RPN are merged with convolutional layers in the middle of the architecture. Yun *et al.* [67] proposed inter- and intra-weighted cross-fusion networks (Infusion-Net), which use a High-Frequency Assistant (HFA) to integrate color and thermal features regarding the feature level gradually. In this procedure, the HFA block exchanges, purifies, and reinforces the object detection-relevant features based on Discrete Cosine Transform (DCT) and Residual Channel Attention Block (RCAB). Additionally, learnable inter- and intra-weight parameters provide optimal information utilization and feature reinforcement for each stream considering each fusion stage. Bao *et al.* [68] proposed a dual-YOLO method based on YOLO v7 [69] for integration of IR and visible images. They also designed attention fusion and fusion shuffle modules to alleviate the false detection rate caused by redundant feature information during the fusion process.

Numerous anchor-free pipelines have recently been proposed for multi-spectral pedestrian detection, which speeds up model detection while avoiding the complex hyper-parameter settings of anchor boxes. Two feature fusion schemes based on a dual-branch *CenterNet* [70] anchor-free detector proposed by [71] for multi-spectral and multi-scale pedestrian detection. The first one is Scale-aware Permutated Attention (SPA) module, which combines local attention and global attention sub-modules, enhancing the quality of the feature fusion at different scales. The second is Adjacent Feature Aggregation (AFA),

which aggregates features across different scales, considering spatial resolution and semantic context. Likewise, Cao *et al.* [72] attempted to train a multi-spectral pedestrian detector without anchor boxes via a box-level segmentation supervised learning framework and compute heat maps. Consequently, the network can be able to localize the pedestrians on small-size input images.

In an innovative approach, Tang *et al.* [73] took illumination into account and designed a progressive image fusion network referred to as *PIAFusion*, which can adaptively maintain the intensity distribution of salient targets as well as retain texture information in the background. It uses an illumination-aware sub-network to estimate the illumination situations and exploits the illumination probability to construct illumination-aware loss. Afterward, the Cross-Modality Differential Aware Fusion (CMDAF) module and halfway fusion strategy merge meaningful information of IR and visible images under the guidance of illumination-aware loss. Likewise, Roszyk *et al.* [74] applied YOLO v4 framework for fast and low-latency multi-spectral pedestrian detection in autonomous driving. Different fusion schemes, as well as different types of models, were investigated, among which feature-level fusion, namely YOLO v4-Middle, demonstrates the best trade-off between accuracy and speed. Peng *et al.* [75] introduced Hierarchical Attentive Fusion Network (HAFNet) embedded with a Hierarchical Content-dependent Attentive Fusion (HCAF) module and a Multi-modality Feature Alignment (MFA) block to overcome the background noise and modality misalignment issue. The MFA exploits the correlation between the TIR and visible domains to fine-tune the pixel alignment of multi-spectral image pairs. Then, the HCAF utilizes top-level features to guide pixel-wise fusion across two streams, resulting in high-quality feature representation. Yadav *et al.* [76] built two uni-modal encoded-decoder feature networks for color and thermal individually using Faster Region-based CNN (Faster R-CNN) [77]. Further, they constructed middle-level CNN fusion architecture, which fused the extracted features in the last convolution layer before feeding it to the decoder for providing the final predictions. Zhang *et al.* [78] presented a Cross-Modality Interactive Attention Network (CIAN) to encode the correlations between two color and thermal spectrums and predict the positions and sizes of pedestrians on a contextual enhanced feature hierarchy. Regarding the halfway fusion strategy, CIAN has investigated three types of operations (*i.e.*, *Elementwise Sum*, *Elementwise Maximization*, and *Concatenation and Channel Reduction*) for how to fuse feature maps. The fusion operation of the *Concatenation and Channel Reduction* shows better performance, which first concatenates the two feature maps, then applies a NIN to reduce the number of channels. To improve the detection accuracy in cases such as occluded objects, light changes, and cluttered backgrounds, Hu *et al.* [79] proposed a Dual-modal Multi-scale Feature Fusion Network (DMFFNet). In their work, MobileNetv3 [80] extracts multi-scale features of dual-modal images as input for MFA module, which processes the spatial information of input feature maps with different scales and establishes longer-distance channel dependencies, thereby reducing background noise interference. Eventually, the Double Deep Feature Fu-

sion (DDFF) module deeply combines the multi-scale features to maximize the correlation between the multi-scale features, which significantly enhances the representation of semantic information and geometric detail.

Cao *et al.* [81] modeled a Multi-spectral Channel Feature Fusion (MCFF) module based on YOLO v4 to fuse the multi-spectral features according to the different illumination conditions. The MCFF module first concatenates the features from visible and thermal modalities in the channel dimension, then uses learning weights to adapt aggregate features.

Gated Fusion Units (GFU) [82] adjusts the contribution of the feature maps generated by each modality via the gating weighting mechanism. Instead of stacking selected features from each channel and adjusting their weights, and motivated by [82], Gated Fusion Double SSD (GFD-SSD) [83] developed two variations of GFU (*i.e.*, *gated fusion* and *mixed fusion*) to fuse the feature maps generated by the two Single Shot MultiBox Detector (SSD) [84] middle layers for multi-spectral pedestrian detection. Using GFUs on the feature pyramid structure, the authors also designed four mixed architectures of both stack fusion and gated fusion (*i.e.*, *Mixed Even*, *Mixed Odd*, *Mixed Early*, and *Mixed Late*), depending on which layers are selected to use the GFUs. By comparing the experimental results on the KAIST dataset, both the GFD-SSD and *Mixed Early* models are superior to the stack fusion. Redundant Information Suppression Network (RISNet) [85] designed a mutual information minimization module to alleviate the influence of cross-modality redundant information on the fusion of RGB-Infrared complementary information. Besides, the RISNet introduced a classification method of illumination conditions based on histogram statistics.

Late fusion-based methods: also known as decision-level fusion, is the high-level fusion technique in which the concatenation is conducted after the last convolutional layer and before fully connected layers or merged the outputs of the two sub-networks such as the location and category prediction. As a work in this category, MultiSpectral Pedestrian DETection TRansformer (MS-DETR) is introduced by [86], which extracts multi-scale feature maps through two parallel modality-specific CNN backbones, aggregates them within the corresponding modality-specific transformer encoders, and fuses the features using a multi-modal transformer decoder. It also adopts a modality-balanced optimization strategy to measure further and balance the contribution of each modality at the instance level. Khalid *et al.* [87] proposed two fusion methods to detect people: In the first one, an encoder-decoder architecture was used for image-level fusion, which independently encodes visible and thermal frames and fed the combined frames into a decoder to produce a single fused image, inputting a Residual Network-152 (ResNet-152) architecture. The second one takes ResNet-152 for feature-level fusion, which extracts features of visible and thermal images separately and concatenates them into a single feature vector as the input of the dense layer. Montenegro *et al.* [88] customized YOLO v5's architecture for low-light pedestrian detection, and also they have conducted experiments on multiple multi-spectral pedestrian datasets *e.g.*, CVC-09, LSI-FIR, FLIR, CVC-14, NightOwls, and KAIST. By extensive evaluations made on different datasets, the best

mean Average Precision (mAP) was obtained on LSI-FIR, followed by CVC-09 and CVC-14. Song *et al.* [89] designed a MultiSpectral Feature Fusion Network (MSFFN) that uses the extracted features of visible-channel and infrared-channel to obtain integrated features. The MSFFN strikes a favorable trade-off between accuracy and speed, especially on small-size input images. They extract multi-scale semantic features using two sub-networks, including Multiscale Feature Extraction of Visible images (MFEV) and Multiscale Feature Extraction of Infrared images (MFEI) and integrate them with an improved YOLO v3 framework. Selective Kernel Network (SKNet) [90] suggested a dynamic selection scheme to adaptively adjust receptive field size using selective kernel units with different kernel sizes. It uses a NIN-based fusion strategy to fuse RGB-IR image pairs. Park *et al.* [91] considered all detection probabilities from RGB, IR, and RGB-IR fusion channels in a unified three-branch model and designed a Channel Weighting Fusion (CWF) and an Accumulated Probability Fusion (APF) layers to fuse probabilities from different information streams at a proposal-level. A combination of an adaptive weight adjustment method with the YOLO v4 [92] is introduced by [93] to enrich the multi-spectral complementarity information for score fusion. Authors in [94] introduced Probabilistic Ensembling (ProbEn), a simple non-learning technique for late-fusion of multiple modalities derived from Bayes' first principle, *i.e.*, conditional independence assumptions. Shaikh *et al.* [95] introduced a probabilistic decision-level fusion approach based on Naïve Bayes to address lighting and temperature changes in color and thermal images by fusion and modeling the detection results of available pedestrian detectors without requiring retraining. In particular, the use of Naïve Bayes for the late fusion strategy enables the network to work with non-registered image pairs as well as poorly registered image pairs.

Zhuang *et al.* [96] examined the impacts of environmental variables on the efficiency of the pedestrian detector and proposed a lightweight Illumination and Temperature-aware Multispectral Network (IT-MN). The method is built on the SSD architecture with designing a late-fusion strategy and Fusion Weight Network (FWN) to compute the fusion weights. In addition, the default box generation is optimized by reducing the number of bounding boxes and choosing specific box aspect ratios to minimize the inference time. Inspired by YOLO v4, Double-Stream Multispectral Network (DSMN) [97] was designed to carry out pedestrian detection in challenging situations such as insufficient and confusing lighting. Their method extracts multi-spectral information provided by RGB and thermal images via two YOLO-based sub-networks. Also, it has an improved Illumination-Aware Network (i-IAN) module to estimate the lighting intensity of varied scenarios and allocate fusion weights to RGB-thermal sub-networks. Li *et al.* [98] explored various fusion schemes and pointed out their key adaptations. They also designed an Illumination-Aware Faster R-CNN (IAF R-CNN) framework to estimate the illumination value of the input image and incorporate color and thermal sub-networks via a gate function defined over the illumination value. In another work, Li *et al.* [99] introduced an Adaptive Soft-Gated Light Perception Fusion (ASG-LPF)

to improve detection performance in varying lighting conditions, which uses a light perception module to distinguish the illumination levels in diverse driving scenarios. Takumi *et al.* [32] proposed a multi-spectral ensemble method based on YOLO v1 [100], which integrates detection results of the four single-spectral detection models into a single space as the final detection. As another approach, LG-FAPF [101] performed a cross-modal feature aggregation process guided by locality information to learn human-related multi-spectral features and used the obtained spatial locality maps of pedestrians as pixel-wise prediction confidence scores for the adaptive fusion of detection results under complex illumination conditions. Considering Center and Scale Prediction Network (CSPNet) model, Wolpert *et al.* [102] proposed an anchor-free multi-spectral framework to investigate various fusion strategies. They also introduced a new data augmentation technique for multi-spectral images called *Random Masking*.

2) **Knowledge Transfer-based Methodologies:** This section's approaches leverage insights from various domains based on knowledge acquired from diverse sources to facilitate nighttime pedestrian detection capabilities. It contains various categories with different methodologies, including transfer learning, supervised and unsupervised domain adaptation, knowledge distillation, and memory-network methods.

Transfer learning methods: transfer learning is reusing the knowledge obtained from pre-trained models for dissimilar but related tasks. In the context of nighttime pedestrian detection, and to fill in the gap of large-scale TIR dataset, Hu *et al.* [103] applied CycleGAN [104] to generate synthetic IR images from visible ones to expand the CVC-09 dataset. They performed experiments using the YOLO v3 and Faster R-CNN models on the CVC-09 dataset, in which the Faster R-CNN has shown better performance in the transfer learning task. In another work by Vandersteegen *et al.* [105], a pre-trained YOLO v2 [106] was used to perform real-time visible-thermal pedestrian detection. Their method takes three image channels composed of a combination of four image channels (*i.e.*, RGB and T) information as input and can work on 80 fps. They discussed the possibility of creating a number of channel combinations as input channels of the YOLO v2 model and designed three models named YOLO-TGB, YOLO-RTB and YOLO-RGT. The YOLO-TGB, which only uses the combination of thermal, green, and blue image channels as input, performs better on the KAIST dataset than other proposed models. Geng *et al.* [107] replaced the loss function of YOLO v3 model with *DIIOU Loss* [108] to accelerate the convergence speed of the network in IR image-based pedestrian detection. Although the loss function curve is more stable, the AP of the Diou-YOLO v3 is not satisfactory.

Domain adaptation methods: the critical idea of employing the domain adaptation mechanism in multi-spectral pedestrian detection is to exploit learned knowledge acquired from the color domain in the thermal images. In this regard, Guo *et al.* [109] focused on image-level domain adaptation by using an image-to-image transformer as a data augmentation tool to convert color images to the thermal spectrum. To aid the joint training process of the domain adapter and the detector, the authors defined a detection loss that back-propagates its

gradients to the image transformer to progressively refine synthetic thermal images. The proposed method provides promising results compared to the baseline on the KAIST benchmark. Kieu *et al.* [110] introduced a task-conditioned training method to help domain adaptation of YOLO v3 to the thermal spectrum. The primary detection network was augmented by adding an auxiliary classification task of day and nighttime thermal images. Additionally, learned representations of this auxiliary task were used to condition YOLO to perform better in the thermal imagery. Authors in [111] addressed three top-down and one bottom-up domain adaptation techniques for pedestrian detection in the nighttime thermal images. They showed that bottom-up domain adaptation achieves better results in challenging illumination conditions. As another work by the same authors [112], a new bottom-up domain adaptation strategy, known as *layer-wise domain adaptation*, is introduced. The main idea for this method is to adjust the RGB-trained detector to adapt to the thermal spectra gradually. Kristo *et al.* [113] attempted to improve the typical object detector performance for person detection at night in challenging weather conditions such as heavy rain, clear weather, and fog. The authors retrained YOLO v3, SSD, Faster R-CNN, and Cascade R-CNN detectors on a dataset of thermal images. They found that YOLO v3 is significantly faster than the others with a processing speed of 27,5 fps. The generalization ability of RPN has been analyzed by [114] for multi-spectral person detection by performing cross-dataset evaluations on several benchmark datasets such as Caltech [115], CityPersons [116], CVC-09, KAIST, OSU, and Tokyo semantic segmentation [117]. They showed that KAIST achieves better results in generalization tasks in both daytime and nighttime conditions.

Unsupervised domain-adaptation methods: the objective of unsupervised domain adaptation is to adapt the well-trained detectors on annotated visible images to the thermal target without any manual annotation effort. As a work in this category, Meta-UDA [118] performed Unsupervised Domain Adaptation (UDA) thermal target detection using an online meta-learning strategy, resulting in a short and tractable computational graph. To mitigate the domain shift between the source and target domain, the Meta-UDA uses the adversarial feature alignment at both the image and instance levels, leading to slight improvement. In another work, Lyu *et al.* [119] used an iterative process to automatically generate the pseudo-training labels from visible and thermal modalities using two single-modality auxiliary detectors. They used the illumination knowledge of daytime and nighttime to assign the fusion priorities of labels for *label fusion*. Without using any manual training labels on the target dataset, the proposed method shows reasonable results on the night scenes of the KAIST dataset. Authors in [120] used transformers to tackle unlabeled data challenges in TIR images. They designed a Self-Supervised Thermal Network (SSTN) to learn feature representation and maximize the mutual data between visible and IR domains by contrastive learning to compensate for the shortage of labeled data. Later, a multi-scale encoder-decoder transformer system was employed for thermal object detection based on the learned feature representations.

Inspired by pseudo-training labels, Lyu *et al.* [121] proposed an unsupervised transfer learning framework in multi-spectral pedestrian detection. Their overall framework is based on a two-step domain adaptation solution, in which the first stage generates intermediate representations of color and thermal images to reduce the domain gap across the source and target domains. The pseudo labels of the target objects are fused via an illumination-aware label fusion mechanism. In the second stage, an iterative fine-tuning process is conducted to progressively converge the detector on the target domain. In another work, Cao *et al.* [122] introduced an auto-annotation framework to iteratively label pedestrian instances in visible and thermal image channels by leveraging the complementary information of multi-modal data. They aim to automatically adapt a pedestrian detector pre-trained on the visible domain to a new multi-spectral domain without manual annotation. The predicted pedestrian labels on both image channels are merged via a label fusion scheme to generate the final multi-spectral pedestrian annotations. Then, the automatically generated labels are fed to a Two Stream Region Proposal Network (TS-RPN) detector to achieve unsupervised learning of complementary semantic features. An unsupervised multi-spectral domain adaptation framework was proposed by Guan *et al.* [123] to generate pseudo-annotations in the source domain, which can be utilized to update the parameters of the model in the target domain according to the complementary information in aligned visible-IR image pairs. Transfer knowledge from thermal to visible domain in unpaired settings and without requiring additional annotations has been performed in [124] by applying image-level and instance-level alignments based on the Faster R-CNN network using adversarial training.

Knowledge distillation methods: The concept of the Knowledge Distillation (KD) is based on inheriting the knowledge learned from a large and complex pre-trained teacher model to a smaller and simpler student model through a supervised learning process [125], [126], [127]. Generally, the main objective of this method is to transfer the applicable and meaningful representations of data to speed up the inference time of the student model without a significant drop in accuracy [126]. According to the teacher-student scheme, Liu *et al.* [128] developed a knowledge distillation framework as a student network that only takes color images as input and generates distinguishing multi-spectral representations, guided by a two-modalities teacher network. Moreover, Cross-modal Feature Learning (CFL) module based on a split-and-aggregation approach was incorporated into the teacher network to learn the standard and modality-specific characteristics between color and thermal image pairs. Hnewa *et al.* [126] employed Cross Modality Knowledge Distillation (CMKD) to enhance the performance of RGB-based pedestrian detection under adverse weather and low-light conditions. Two different CMKD methods were developed to transfer the multi-modal information of a teacher detector to a student RGB-only detector. The former uses KD loss, while the latter integrates adversarial training with knowledge distillation. Zhang *et al.* [127] proposed a Modality Distillation (MD) framework to transfer the knowledge from a high thermal resolution two-stream network with feature-level fusion to a low thermal

resolution single-stream network with early fusion strategy. In particular, two specific knowledge distillation modules are used in the MD framework. An attention transfer generates attention masks by GAFF from a two-stream teacher model, which is transferred to a single-stream student model through performing an early fusion. Finally, a semantic transfer resolves the problem of modality imbalance in feature distillation using a new Focal Mean Square Error (F-MSE) cost function.

Memory-network methods: Memory Augmented Neural Network (MANN) can memorize and recall the prior information, such as visual appearance in the memory module, so the relevant data can be accessed by calculating the similarity [129]. In [130], a pedestrian detection process is introduced to improve the detector’s performance in any modality. In the first stage, a multisensory-matching contrastive loss guides the pedestrian visual representation of two visible and thermal modalities to be similar. In the second, a Multi-Spectral Recalling (MSR) memory improves the visual representation of the single modality features by recalling the visual appearance of multi-spectral modalities and memorizes the multi-spectral contexts through a multi-spectral recalling loss, which encoded more discriminative information from a single input modality. The Large-scale Pedestrian Recalling (LPR) based on key-value memory was proposed by [129], which memorizes visual information of large-scale pedestrians to recall the relevant characteristics to cover inadequate small-scale pedestrian appearances.

3) **Image Specification-based Methodologies:** Another category focuses on methodologies in which image specifications play a crucial role. These methods can be divided into three primary strategies: *image enhancement*, *image-to-image translation*, and *saliency maps* methods.

Image enhancement methods: TIR images are characterized by noisy details, blurred edges, low contrast, and low resolution, resulting in a performance drop caused by low discrimination. In this regard, low-light image enhancement techniques are considered to improve the visual quality of thermal images and simplify their challenges. In a work by Marnissi *et al.* [131], an enhancement method based on images’ architecture, title Thermal-Enhancement GAN (TE-GAN) is designed, which constituted of contrast augmentation, noise elimination, and edge restoration. To enhance the clarity of the IR pedestrian targets with blurred edges, Sun *et al.* [132] adopted a super-resolution algorithm called Wide Activation Deep Super-Resolution (WDSR)-B [133]. They add the four-time down-sampling layer output to YOLO v3 trained by the enhanced IR images to acquire richer context information for small pedestrian targets. In another work, Marnissi *et al.* [134] combined Generative Adversarial Network (GAN) and Vision Transformer (ViT) for thermal image enhancement and introduced Thermal Enhancement Vision Generative Adversarial Network (TE-VGAN). TE-VGAN employs the U-Net architecture as an input image generator and two ViT models as global and local discriminators. The thermal loss feature is also introduced in their work to generate high-quality images. They investigated the effect of the thermal image enhancement method on the detection performance of different YOLO

versions, resulting in a balance between contrast enhancement and noise reduction.

DIVFusion [135] incorporates a low-light image enhancement task and a dual-modal fusion task in a unified framework to investigate the effect of lighting conditions on image fusion. In their method, firstly, a Scene-Illumination Disentangled Network (SIDNet) is devised to eliminate the illumination degradation in nighttime visible images while maintaining informative features of source images. Then, a Texture-Contrast Enhancement Fusion Network (TCEFNet) is employed to aggregate complementary information and boost fused features’ contrast and texture details. Finally, a color consistency loss is used to alleviate color distortion in enhancement and fusion processes. Li *et al.* [136] built Feature Attention Module (FAM) and Feature Transformation Module (FTM) to improve the efficiency of a pedestrian detector in darkness. FAM is designed to suppress the noisy representations, while FTM allows pedestrian examples under a low-light environment to generate more discriminate feature representations. An attention-based feature fusion module was designed in [137] to enhance pedestrian detection in low-illumination images. They used the brightness channel (*i.e.*, V-channel) from the *HSV* image of the thermal image as an attention map to activate the unsupervised auto-encoder for obtaining more details about the pedestrian. In order to address the challenge of light compensation in low-light conditions, a Brightness Correction Processing (BCP) algorithm is considered to guide self-attention map learning. Eventually, the image enhancement method was integrated into YOLO v4 detection model. They evaluated the proposed architecture on the LLVIP dataset.

To highlight pedestrians in low-resolution and noisy IR images, an Attention-guided Encoder-Decoder Convolutional Neural Network (AED-CNN) [138] is devised. In AED-CNN, the encoder-decoder module generates multi-scale features, and a skip connection block is integrated into the decoder to fuse the feature maps from the encoder and decoder structure. By adding an attention module, the network effectively emphasizes informative features and suppresses background interference while re-weighting the multi-scale features generated by the encoder-decoder module. Patel *et al.* [139] introduced a computationally compact algorithm based on Depthwise Convolution (DC) with the aim of network parameters reduction. The proposed algorithm enhances the details of the thermal images using Adaptive Histogram Equalization (AHE) and extracts the salient features in these images by a new Convolutional Backbone Network (CBN), where depthwise convolution minimizes the computational complexity. YOLO-FIRI [140] is another method developed for pedestrian detection in IR images, which achieved outstanding results by making improvements on YOLO v5 structure. Firstly, by extending shallow CSPNet in the backbone network and incorporating an improved Select Kernel (SK) attention module in the residual block, it forces the model to focus on shallow and detailed information and learn the distinguishable features. Secondly, the detection accuracy of small and blurry pedestrians in IR images is increased by adding four-scale feature maps to the detection head. Finally, Densefuse [141] is adopted as a data enhancement to fuse visible and infrared images to boost the

features of IR images.

Image-to-image translation methods: the goal of Image-to-Image (I2I) translation models is to learn the visual mapping between a source and target domain while preserving the essential features. Specifically, I2I has been widely used in image colorization, denoising, and synthesis [142]. In these approaches, thermal image colorization aims to translate from the temperature-channel domain into the RGB channel. PearlGAN presented in [143] to facilitate the translation of nighttime Thermal-Infrared (TIR) image into a daytime color one. By taking advantage of a top-down guided attention module and a corresponding attention loss, *PearlGAN* can produce hierarchical attention distribution and reduce local semantic ambiguity in IR images through context information. In addition, a structured gradient alignment loss was designed to enhance edge consistency during the translation. The colorization of thermal-IR images in pedestrian detection application is accomplished by [144], organized into three main modules: thermal image colorization, improvement of colorized images, and pedestrian detection. The colorized and improved images are fed to the detection head using a pre-trained YOLO v5 framework.

To mitigate color distortion and edge blurring caused by translation from temperature spectrum to color spectrum, [145] considered a one-to-one mapping relationship and introduced an improved CycleGAN [104], called Gray Mask Attention-CycleGAN (GMA-CycleGAN). It first translates the TIR images to Grayscale Visible (GV) and then uses the original CycleGAN to obtain the translation from GV to Color Visible (CV). A mask attention module based on the thermal temperature mask and the color semantic mask has been designed without increasing training parameters to better differentiate between pedestrians and the background. Meanwhile, to make the texture and color of the translated image more realistic in the feature space, a perceptual loss was added to the original CycleGAN loss function. Devaguptapu *et al.* [146] proposed to borrow knowledge from the large-scale RGB dataset without the need for paired multi-modal training examples and used CycleGAN to implement an unpaired image-to-image translation framework. It can generate pseudo-RGB equivalents of a given thermal image and employs a multi-modal Faster R-CNN detector for pedestrian detection in thermal imagery. To transform the visible domain into the thermal domain, authors in [147] implemented a generative data augmentation method based on the Least-Squares GAN (LS-GAN) [148]. They also used the perceptual loss function to measure the similarity between authentic and synthesized images in pixel space.

Saliency maps methods: the purpose of salient object detection is to highlight the most noticeable areas in the given image and distinct the prominent objects from their surroundings using the intensity of each pixel. Accordingly, in TIR images, the saliency maps can be used to detect temperature. Altay *et al.* [149] presented a two-branch architecture that can incorporate features of thermal images with their correlated saliency maps to acquire better representations of pedestrian regions. Instead of using color-thermal image pairs in the fusion network, Ghose *et al.* [150] suggested augmenting thermal images with their corresponding saliency maps, which

produced by static methods and two deep saliency networks, Pixel-wise Contextual Attention Network (PiCANet) [151] and Recurrent Residual Refinement Network (RRRNet) [152]. Marnissi *et al.* [153] proposed a bi-spectral image fusion scheme, which was augmented with a corresponding saliency map using Visual Salient Transformer (VST) and also incorporated this fusion process into the YOLO v3 as base architecture for real-time applications. The proposed approach has shown its advantage in low computational cost, which allows faster inference time. Zhao *et al.* [154] put more emphasis on the temperature information in infrared images by constructing an IR-temperature transformation formula which can convert the IR images into corresponding temperature maps. It finally uses a trained temperature network for pedestrian detection. On the OSU and FLIR datasets, the transformed temperature maps boost the overall performance regardless of external influences.

4) **Images Discrepancy-based Methodologies:** These methods target enhancing the accuracy and reliability of nighttime pedestrian detection by exploiting discrepancies within images and characteristics of different imaging sensors and analyzing variations in image quality and content. The modality discrepancy is alleviated by focusing on *modality imbalance problem* and *position-shift problem*.

Modality imbalance problem: Scenes in which one sensor performs considerably better than the others can lead to a *bias* in training towards one dominant input modality. For instance, uneven distribution of training data in multi-modal learning causes less contribution of the non-dominant input modality during network training and, therefore, limits the generalizability of the model. In this regard, Oksuz *et al.* [155] provided a comprehensive taxonomy of the imbalance problems in object detection. They categorize these problems into four significant categories: *class imbalance* (i.e., inequality distribution of training data among different classes), *scale imbalance* (i.e., various scales of objects), *spatial imbalance* (i.e., spatial properties of the bounding boxes), and *objective imbalance* (i.e., minimization of multiple loss functions). Additionally, in multi-spectral pedestrian detection, the modality imbalance issue substantially impacts the algorithm performance, which can occur in two different ways, including the *illumination modality imbalance problem* and the *feature modality imbalance problem* [156]. Das *et al.* [157] proposed a training process with a regularization term i.e., *Logarithmic Sobolev Inequalities* [158] to consider the features of both modalities equally during fusion. The proposed regularizer reduces the modality imbalance in the network by equally distributing the training data among the modalities. Li [159] trained YOLO v3 framework to detect pedestrians under insufficient illumination conditions. In their method, focal loss [64] was added to the loss function to overcome the imbalance issue of IR images. Zhou *et al.* [156] resolved the modality imbalance issue in multi-spectral images through the implementation of a single-stage Modality Balance Network (MB-Net), which included a Differential Modality Aware Fusion (DMAF) and an Illumination Aware Feature Alignment (IAFA) module to extract complementary information and align the two modality features according to the lighting conditions. Dasgupta *et al.*

[160] developed Multi-modal Feature Embedding (MuFEM) module using Graph-Attention Network (GAT) [161] to deal with the imbalance issue between color image branch and thermal image branch. Also, the channel-wise attention block and four-directional IRNN (4Dir-IRNN) block [162] are incorporated in Spatio-Contextual Feature Aggregation (SCoFA) to improve fusion using spatial and contextual information of the pedestrian. The 4Dir-IRNN block consists of four Recurrent Neural Networks (RNNs), which compute context features in four directions.

Position-shift problem: The physical properties of different cameras (e.g., Field-of-View (FoV), resolutions, wavelengths, etc.) can cause weakly aligned image pairs in multi-spectral data, where the positions of the objects are out of synchronization on different modalities. Some works tried to address the mentioned problem in multi-modal sensors using geometrical calibration and image alignment methods. The study by Zhang *et al.* [163] is the first work providing insights into the position shift problem between color and thermal images. They introduced an Aligned Region CNN (AR-CNN) detection framework to solve the weakly aligned image pairs. The AR-CNN firstly predicts the position shift and adaptively aligns the region feature maps of the two modalities through a Region Feature Alignment (RFA) module. Based on the aligned features, a confidence-aware fusion method is proposed to accomplish feature re-weighting, which selects the highly informative features while suppressing the useless ones. Moreover, a ROI jitter strategy is adopted to enhance the robustness of position shift patterns. Kim *et al.* [164] used adversarial learning to make each spectrum share its complementary information in a common feature space to compensate for the lack of aligned multi-spectral pedestrian datasets. Kim *et al.* [165] have constructed uncertainty-aware multi-spectral pedestrian detection architecture to handle miscalibration (i.e., different FoV in color and thermal cameras) and modality discrepancy challenges. For the miscalibration issue, the Uncertainty-aware Feature Fusion (UFF) module was formulated to mitigate the impact of ambiguous Region of Interest (ROI). The modality discrepancy is alleviated through the Uncertainty-aware Cross-modal Guiding (UCG) module, which can encode more discriminative visual representations. Wanchaitanawong *et al.* [166] introduced a multi-modal Faster R-CNN robustly against significant misalignment between the two modalities. The key points are modal-wise regression for bounding-box regression of each modality to deal with the significant misalignment and multi-modal Intersection over Union (IoU) for mini-batch sampling that combines the IoU for both modalities.

5) **Multi-task methods:** Multi-task learning is a training paradigm that aims to learn multiple related tasks simultaneously, using shared feature representations [167]. A cross-task feature alignment method was proposed by [168] to tackle the misalignment of scale and channel of features from image relighting and pedestrian detection tasks by placing four feature alignment layers before the feature fusing and sharing step in cross-task learning. Meanwhile, a multi-scale feature-enhanced detection network expands the receptive field of the multi-scale feature extractor and thereby provides richer semantic information of fused features for the detection head.

An illumination-aware weighting mechanism is presented by Guan *et al.* [169] to adaptively re-weight the detection results of day- and night-illumination sub-networks to learn multi-spectral human-related characteristics to perform pedestrian detection and semantic segmentation under various illumination conditions, simultaneously. Dai *et al.* [170] developed the Faster R-CNN detector using the ResNet-50 as a feature extractor for pedestrian detection and distance estimation using the NIR-based camera. An Automatic Region Proposal Network (ARPN) was designed by [171] to get bounding boxes. A pedestrian segmentation task is also added based on a Feature Pyramid Network (FPN) [172] to obtain the confidence scores. To distinguish pedestrian examples from complex negative samples, Li *et al.* [173] added two sub-networks for jointly semantic segmentation and pedestrian detection tasks to the unified fusion network, which is denoted as *MSDS-RCNN*. The paper also studied the effects of training annotation noise by creating a sanitized version of KAIST ground-truth annotations so that the sanitized training annotations significantly reduce the inference error. Evaluations showed that the segmentation supervision benefits multi-spectral pedestrian detection.

6) **Other methods:** As for the final category, this subsection introduces works that cannot fit into the previous ones. Accordingly, the authors in [174] employed a region decomposition branch in Faster R-CNN architecture, which exploits the multi-region features, including head, body trunk, and legs, to solve the pedestrian occlusion problem in thermal images. The proposed architecture learns the high-level semantic features by combining the global and partial appearance features step by step. The Center and Scale Prediction Network (CSPNet) [175] has been applied in [176] to obtain three IR pedestrian detection models, namely daytime, nighttime, and full-time. The full-time model has a lower detection loss rate, while the nighttime model and the daytime model perform poorly in detecting small objects in the evening, respectively. Xu *et al.* [177] aggregated ground-area context information into the Faster R-CNN for pedestrian detection and shared the predicted ground horizon area to a Ground-Region Proposal Network (GRPN), which can only process the pixels on the proposed horizon region to minimize False Positive (FP) rate. Since the output of the *FC layer* is the position vector of pixels in the horizon region, the size of the GRPN model is largely increased and has a high computational cost. Dai *et al.* [178] compared and analyzed visible and IR images acquired by using visible-spectrum, Near-Infrared (NIR), Short-wave Infrared (SWIR), and Long-wave Infrared (LWIR) cameras. For the first time, they used a nine-layer CNN model with a self-learning SoftMax [179] to detect nighttime pedestrian samples in NIR images. In order to enhance the detection accuracy of multi-scale pedestrians in IR images, in [180], two regional proposal networks based on the Faster R-CNN architecture were designed to focus on near and far away pedestrians. Although the proposed multi-scale RPN has shown improvements in far-away pedestrian detection, it is not optimized to work in real time. Kalita *et al.* [181] have presented a real-time human detection system using YOLO v3, which achieved a speed of 17 millisecond per

image on the KAIST thermal dataset. The brightness aware Faster Region-based CNN (Faster R-CNN) model [182] was proposed to perform the pedestrian prediction under low-light and day-light scenarios. In the first step, the model calculates the brightness of the input image based on the pixel intensity to predict the day or night scenario. In the second step, two separate thermal or color models are employed for pedestrian detection based on the first step output. It should be noted that the authors trained the FLIR dataset for the thermal model and the PASCAL VOC dataset for the color model.

C. Hybrid Approaches

Hybrid methods combine elements of handcrafted features and deep learning, aiming to harness the strengths of each approach for improved nighttime pedestrian detection performance. Accordingly, they require significantly less computational resources than deep learning methods and overcome the poor generalization of handcrafted methods. However, the performance of such approaches is not properly optimized in terms of prediction accuracy or model running time.

As a hybrid methodology for nighttime pedestrian detection, the study of Kim *et al.* [183] presented a method to detect pedestrians at night using a visible-light camera and Faster R-CNN model, which can handle the changes of the pedestrians' spatial position by fusing deep convolutional features in successive frames. To make the model robust against noise and illumination, the authors used Additive Random White Gaussian Noise (AWGN) and applied two pre-processing methods, *i.e.*, pixel normalization and Histogram Equalization (HE) mean subtraction, to normalize the illumination and contrast levels of successive frames. Besides, a weighted summation of successive frame features was added to exploit temporal information about the pedestrian, which enhanced the accuracy of the pedestrian detector at nighttime. To find the optimal fusion stage in CNN, authors in [184] used RPN to merge the features of visual and IR images. After halfway feature fusion in RPN, they employed Boosted Decision Tree (BDT) classifier to improve pedestrian detection results and reduce the false positive rate. Tumas *et al.* [185] eliminated the sliding window technique and applied background subtraction to extract thermally active points as Region of Interest (ROI) for pedestrian detection in Far-Infrared (FIR) domain. The proposed technique accelerates the Histogram of Oriented Gradients (HOG) based pedestrian detector to run at 6 fps using only CPU performance. Narayanan *et al.* [186] developed a model for low-light pedestrian prediction using HOG and YOLO v3 algorithm. They also experimented the detection accuracy of HOG detector and SVM classifier in thermal images. Xu *et al.* [187] designed a framework for learning and transferring cross-domain feature representations for pedestrian detection that works based on two different networks. The first one, titled Region Reconstruction Network (RRN), is employed to learn a non-linear feature mapping and model the relations among the color and IR image pairs. Afterward, the cross-modality feature representations learned from RRN are transferred to a second network titled Multi-Scale Detection Network (MSDN), which operates only on

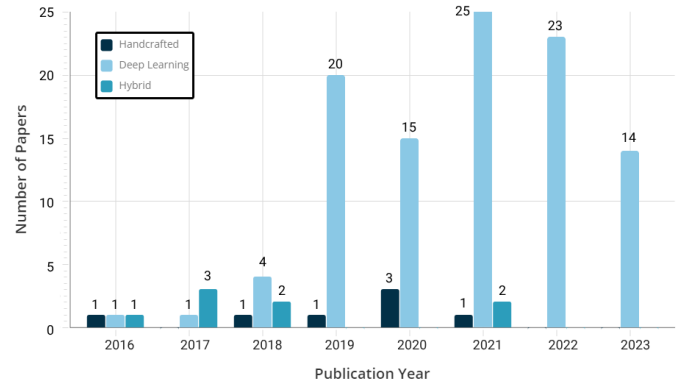


Fig. 7: The trends of employing various approaches explained in this survey by state-of-the-art research works published in different years.

RGB inputs and outputs the recognition results. Both RRN and MSDN networks have employed ACF [188] to generate pedestrian proposals. In this way, only color images are considered at the test phase, and no thermal data are needed, which significantly reduces the cost of thermal data annotation. In [189], Support Vector Regression (SVR) was adopted to learn the pedestrians probabilities, which performs well on small-scale pedestrians. Chen *et al.* [190] utilized a Total Variation (TV) minimization [191] method based on structure transfer to integrate TIR-RGB image pairs, preserving the infrared intensity distribution and the local appearance features. However, when the thermal radiation of the pedestrian and the background are the same, the performance of the detector is affected.

V. DISCUSSION

Regarding the state-of-the-art methodologies introduced in previous sections, this section discusses the current trends and future expectations of the works targeting nighttime pedestrian detection.

A. Employed Methodologies Trends

Regarding categorizing the works introduced in Section IV, nighttime pedestrian detection approaches can be divided into handcrafted features, deep learning, and hybrid methodologies. In this regard, Fig 7 shows the distribution of the surveyed paper regarding the primary categories they belong to. According to the figure, it can be seen that the majority of the works published in the last two years consider deep learning-based techniques the most reliable methodology to detect pedestrians in low-light conditions. In other words, recent approaches only focus on employing DNNs instead of handcrafted and hybrid approaches. The main reason may be attributed to *automatic learning of features* in DNNs, which cover many possible conditions in which pedestrians are challenging to detect. Additionally, the works are getting more practical, providing the possibility to be used in real applications, and making domain-specific applications based on handcrafted or hybrid methods is not a practical solution.

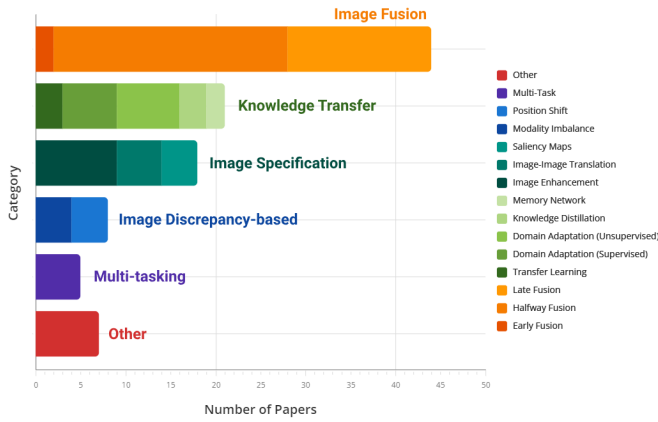


Fig. 8: Distribution of the reviewed papers considering the sub-categories introduced in Section IV.

Moreover, in a more detailed chart, Fig 8 shows the distribution of papers surveyed regarding the introduced sub-categories. It can be seen that most of the papers have targeted *image fusion* techniques for nighttime pedestrian detection applications. Since thermal imaging (*i.e.*, long wavelength IR) can capture the infrared radiation from objects and are sensitive to temperature changes, thermal images provide clearer contours information of pedestrians under insufficient lighting conditions. However, the thermal IR modality lacks visual details such as texture, color, and precise edges of the objects, which can be captured by RGB sensors. In addition, the quality of visible images is significantly degraded under severe weather conditions, low resolution, and unfavorable lighting. Considering the characteristics of both visible and thermal sensors, cross-spectral fusion has become a promising alternative solution for overcoming the limitations of an unimodal approach to adapt to the all-weather and all-day situations. By fusing complementary visual features from multiple modalities, the stability, reliability, and perceptibility of the pedestrian detectors are enhanced. Despite the great progress made in multi-spectral pedestrian detection, there still exists a large gap between the current artificial vision systems and human vision ability. Among them, *halfway fusion* covers most of the works, and *late fusion* is the second preferred approach in *image fusion* subcategory. *Knowledge transfer* and *image specification* methodologies are the following trendy solutions according to the stats in the figure. It can also be seen that *multi-task* methodologies have not absorbed massive attention among the papers published in recent years in the domain.

It should also be added that three types of deep learning-based architectures have been dedicated to achieving multi-spectral pedestrian detection, which can be categorized into the conventional CNN-based, Auto-Encoder (AE)-based, and GAN-based architectures. Fig 9 shows the distribution of these methodologies and architectures in brief. Accordingly, end-to-end CNN-based methods contain feature extraction, feature fusion, and image reconstruction processes through well-designed loss functions and network architectures. On the other hand, the AE-based methods first train the encoder

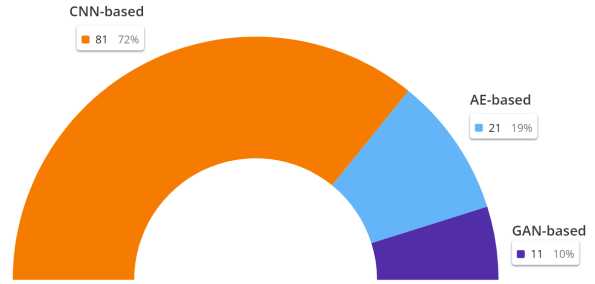


Fig. 9: Distribution of deep learning-based architectures for multi-spectral pedestrian detection.

and the decoder as the feature extractor and the image reconstructor, respectively. Then, the multi-image fusion process is accomplished according to the fusion rules. Finally, and in the GAN-based methods, the architecture is suitable for unsupervised pedestrian detection, relying on the adversarial mechanism between the generator and discriminator. The discriminator forces the generator to make the target distribution in the fused images as close as possible to the source images.

B. Dataset Trends

Regarding the datasets introduced in Section III, the surveyed research works have been evaluated on various datasets. Thus, Fig 10 depicts the distribution of the utilized datasets by the reviewed papers. It can be seen that most of the papers (*i.e.*, around half of them) have utilized KAIST dataset. The second and third datasets other research works use are FLIR and CVC-14, respectively. It should be mentioned that some of the research works (*i.e.*, ~ 5.3 percent) prefer to evaluate their in-house collected, mainly collected from real-world scenarios. As introduced in Section IV, the differences in the physical characteristics of sensors lead to the misalignment of image pairs and have limited applicability in real-life situations. Despite the increasing number of visible-IR datasets in recent years, accessing instances with strictly aligned multi-spectral images is still a challenging problem. The benchmark datasets reported in the scientific literature can only provide information under certain scenes, most of which are recorded by a stationary camera. Therefore, there is a lack of datasets that contain a sufficient variety of fine-grained annotated samples taken from a moving camera, as environments can change dynamically.

C. Performance Evaluations

Considering the performance of the surveyed works, Table II analyzes the computational efficiency of the state-of-the-art methods on KAIST test set. It should be noted that the best and second-best results are boldfaced and underlined, respectively. According to the table, the MD [127] method tops the chart in processing speed and takes only 0.007 seconds to process a single image. The main reason for such performance is due to the theory of knowledge distillation, which accelerates inference by transferring the knowledge learned from a high thermal-resolution model to a low one. The second best result

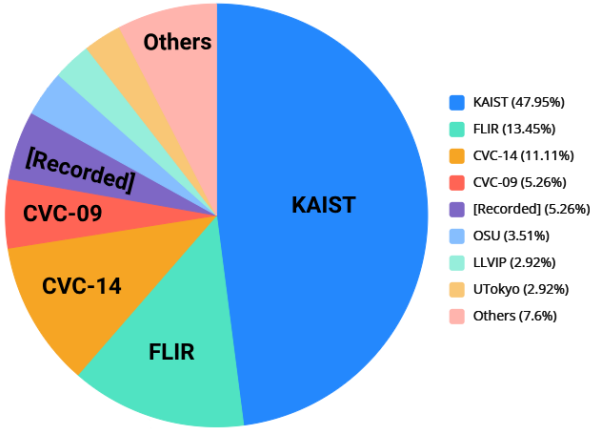


Fig. 10: Distribution of the datasets available for evaluation of nighttime pedestrian detection in different works.

is the GAFF [61] model, which requires only 0.009 seconds of inference time. The main reason for such performance is that the GAFF only includes three convolution layers, so the total number of learnable parameters and the computational cost is low. There are some approaches with performances negligibly less than these approaches, including YOLO-TGB [105] with 0.012, Dual-YOLO [68] with 0.016, HAFNet [75] with 0.017, and Marnissi *et al.* [153] with 0.019 seconds to process a single image. On the other hand, CMT-CNN [187] as a hybrid approach is the most computationally intensive methodology, with 0.59 seconds to process a single image. The main reason for such low performance is due to the use of ACF proposals at the test time, which is a time-consuming process.

Moreover, Table III shows the detection accuracy in evaluating different approaches. The results are reported in terms of MR under *Reasonable* settings, and the approaches are classified according to the categories presented in Section IV. As shown in the Table, the MCFF [81] ranks first as a halfway-fusion strategy in overall performance on the KAIST by a large margin. The main reason for such performance is due to the MCFF transferring the fusion information from the bottom to the top at different stages. It can be observed that in the *Reasonable* nighttime criteria, the MCFF [81] obtains superior results than its daytime experiment. The main reason for such performance is that the MCFF uses the illumination information to learn the fusion weights. Similarly, LG-FAPF [101] as a late-fusion strategy performs remarkably better compared to the other detectors. The main reason for such performance is due to a locality-guided pixel-level fusion scheme that aggregated the human-related features in the complementary modalities to integrate the prediction confidence scores in color and thermal channels. Among these methods, only four methodologies (*i.e.*, CMT-CNN [187], Kim *et al.* [183], Choi *et al.* [189], and Chen *et al.* [190]) are hybrid approaches, which witnessed a significant drop in MR. It can be concluded that the hybrid approaches are not applicable to around-the-clock applications, and specifications are required.

As the final discussion, it is essential to note that by

TABLE II: Performance evaluation of state-of-the-art deep learning-based multi-spectral pedestrian detectors on KAIST test set. The superscripts *X*, *V*, *K*, *P*, *I*, *2*, and *3* represent NVIDIA GPU models used for evaluation, including *TitanX*, *Tesla V100*, *Tesla K40*, *Tesla P40*, *1080Ti*, *2080Ti*, and *3090Ti*, respectively. The best and second-best results are boldfaced and underlined, respectively. Additionally, *s/f* represents seconds per frame.

Family	Backbone	Method	Speed (s/f)
DL	CSPDarknet-53	ASPFF Net [56]	0.028 ¹
	CSPDarknet-53	MCFF [81]	0.031 ^P
	CSPDarknet-53	YOLO-CMN [55]	0.02 ²
	CSPDarknet-53	Chan <i>et al.</i> [93]	0.76 ³
	CSPDarknet-53	DSMN [97]	0.76 ³
	Custom CNN	RISNet [85]	0.1 ^V
	Darknet-19	YOLO-TGB [105]	0.012 ¹
	Darknet-53	TC Det [110]	0.033 ¹
	Darknet-53	Marnissi <i>et al.</i> [153]	0.019 ^X
	ELAN	Dual-YOLO [68]	0.016 ³
	MobileNet v2	IT-MN [96]	0.03 ^X
	MobileNet v3	DMFFNet [79]	0.021 ²
	ResNet-18	MD [127]	0.007¹
	ResNet-50	Zou <i>et al.</i> [71]	0.04 ¹
	ResNet-50	MB-Net [156]	0.07 ¹
	ResNet-50	BAANet [60]	0.07 ¹
	ResNet-50	HAFNet [75]	0.017 ¹
	ResNet-50 + FPN	ProbEn [94]	0.025 ²
	ResNet-101	ResNet + FPN [63]	0.129 ^X
	VGG-16	CIAN [78]	0.066 ¹
	VGG-16	Kim <i>et al.</i> [165]	0.11 ¹
	VGG-16	GFD-SSD [83]	0.0512 ¹
	VGG-16	AR-CNN [163]	0.12 ¹
	VGG-16	GAFF [61]	0.009 ¹
	VGG-16	MSR [130]	0.04 ¹
	VGG-16	Park <i>et al.</i> [91]	0.58 ^X
	VGG-16	Halfway Fusion [58]	0.43 ^X
	VGG-16	IATDNN+IASS [169]	0.25 ^X
	VGG-16	IAF R-CNN [98]	0.21 ^X
	VGG-16	MSDS-RCNN [173]	0.22 ^X
	VGG-16	LG-FAPF [101]	0.14 ^X
	VGG-16	Ding <i>et al.</i> [65]	0.222 ^X
	VGG-16	HMFN [72]	0.026 ^X
	VGG-16	Ding <i>et al.</i> [90]	0.071 ^X
Hybrid	VGG-16+ACF	CMT-CNN [187]	0.59 ^K

expanding the use of fully autonomous vehicles, the challenges of correct and real-time detecting pedestrians under various scenarios are becoming inevitable. Accordingly, explainable and interpretable mechanisms to exploit why a system failed/succeeded in a scenario can bring about more public reliability and confidence among people, including pedestrians, for interacting with autonomous systems. Thus, tailoring the current methodologies with the field of Explainable AI (xAI) is another direction to be investigated by researchers.

VI. CONCLUSIONS

The paper in hand provided a comprehensive survey of pedestrian detection approaches tailored to low-light conditions, addressing a crucial challenge in computer vision, surveillance, and autonomous driving. The accurate and reliable recognition and tracking of pedestrians under reduced visibility is of paramount importance for enhancing the safety of autonomous vehicles and preventing accidents. The survey

TABLE III: Miss Rate (MR) comparison of state-of-the-art deep learning-based multi-spectral pedestrian detectors in three subsets of the KAIST test set, *i.e.*, all-day, day-time, and night-time. The best and second-best results are boldfaced and underlined, respectively. Note that the lower MR is better.

Method	Family	Category	Backbone	All-Day	Day-Time	Night-Time
Halfway Fusion [58]	DL	Halfway-Fusion	VGG-16	36.99	36.84	35.49
Yadav <i>et al.</i> [76]		Halfway-Fusion	VGG-16	29.00	26.00	32.00
GFD-SSD [83]		Halfway-Fusion	VGG-16	28.00	25.80	30.03
CFR [27]		Halfway-Fusion	VGG-16	<u>6.13</u>	7.68	3.19
CIAN [78]		Halfway-Fusion	VGG-16	27.71	30.74	21.07
Ding <i>et al.</i> [65]		Halfway-Fusion	VGG-16	34	36	35
GAFF [61]		Halfway-Fusion	ResNet-18	7.93	9.79	4.33
BAANet [60]		Halfway-Fusion	ResNet-50	7.92	8.37	6.98
CS-RCNN [54]		Halfway-Fusion	ResNet-50	11.43	11.86	8.82
HAFNet [75]		Halfway-Fusion	ResNet-50	6.93	7.68	5.66
Zou <i>et al.</i> [71]		Halfway-Fusion	ResNet-50	7.77	9.41	2.00
ResNet-101 + FPN + Sum [63]		Halfway-Fusion	ResNet-101	27.60	27.92	25.77
Yang <i>et al.</i> [53]		Halfway-Fusion	ResNet-101	10.71	13.09	8.45
YOLO-CMN [55]		Halfway-Fusion	CSPDarknet-53	7.85	8.03	7.82
MCFF [81]		Halfway-Fusion	CSPDarknet-53	4.91	6.23	2.90
ASPPF Net [56]		Halfway-Fusion	CSPDarknet-53	11.64	14.14	6.73
DMFFNet [79]		Halfway-Fusion	MobileNet v3	9.26	12.79	5.17
RISNet [85]		Halfway-Fusion	Custom CNN	7.89	<u>7.61</u>	7.08
IAF R-CNN [98]	DL	Late-Fusion	VGG-16	15.73	14.55	18.26
Ding <i>et al.</i> [90]		Late-Fusion	VGG-16	32	34	34
LG-FAPF [101]		Late-Fusion	VGG-16	5.12	5.83	3.69
Park <i>et al.</i> [91]		Late-Fusion	VGG-16	31.36	31.79	30.82
ProbEn [94]		Late-Fusion	ResNet-50+FPN	7.66	9.07	4.89
MS-DETR [86]		Late-Fusion	ResNet-50+ResNet-18	<u>6.13</u>	<u>7.78</u>	3.18
DSMN [97]		Late-Fusion	CSPDarknet-53	14.33	13.34	22.36
IT-MN [96]		Late-Fusion	MobileNet v2	14.19	14.30	13.98
Das <i>et al.</i> [157]	DL	Modality-Imbalance	PVT	7.41	7.69	7.03
Dasgupta <i>et al.</i> [160]		Modality-Imbalance	ResNeXt-50	9.23	9.33	8.97
MB-Net [156]		Modality-Imbalance	ResNet-50	<u>8.13</u>	<u>8.28</u>	<u>7.86</u>
Kim <i>et al.</i> [164]	DL	Position-Shift	ResNet-50	42.89	42.42	43.65
AR-CNN [163]		Position-Shift	VGG-16	9.34	9.94	8.38
Kim <i>et al.</i> [165]		Position-Shift	VGG-16	8.45	9.39	7.39
Wanchaitanawong <i>et al.</i> [166]		Position-Shift	VGG-16	9.67	10.69	9.24
VGG-16-two-stage [109]	DL	Domain-Adaptation	VGG-16	46.30	53.37	31.63
BU (VLT, T) [112]		Domain-Adaptation	Darknet-53	25.61	32.69	<u>10.87</u>
TC-Det [110]		Domain-Adaptation	Darknet-53	<u>27.11</u>	<u>34.81</u>	10.31
Marnissi <i>et al.</i> [124]	DL	Unsupervised Domain-Adaptation	ResNet-101	44.60	50.29	28.79
U-TS-RPN [122]		Unsupervised Domain-Adaptation	VGG-16	36.42	37.15	33.00
UTL [121]		Unsupervised Domain-Adaptation	VGG-16	19.98	22.17	15.78
Feature-Map Fusion [119]		Unsupervised Domain-Adaptation	VGG-16	<u>23.09</u>	<u>24.55</u>	<u>17.74</u>
Kim <i>et al.</i> [129]	DL	Memory-Network	VGG-16	19.16	24.70	8.26
MSR [130]		Memory-Network	ResNet-101	10.32	13.28	6.23
IATDNN+IASS [169]	DL	Multi-Task	VGG-16	26.37	27.29	24.41
MSDS-RCNN [173]		Multi-Task	VGG-16	11.63	10.60	13.73
DCRL-PDN [128]	DL	Knowledge-Distillation	VGG-16	25.89	27.01	23.82
MD [127]		Knowledge-Distillation	ResNet-18	8.03	9.85	4.84
LS-GAN [147]	DL	I2I-Translation	Darknet-53	25.62	31.86	12.92
YOLO-TGB [105]	DL	Transfer-Learning	Darknet-19	31.2	34.7	23.1
Ghose <i>et al.</i> [150]	DL	Saliency-Maps	VGG-16	-	30.4	21.0
Song <i>et al.</i> [176]	DL	Other	ResNet-50	-	12.23	4.56
CMT-CNN [187]	Hybrid	-	VGG-16+ACF	49.55	47.30	54.78
Choi <i>et al.</i> [189]		-	VGG-16+ACF	47.31	49.31	43.75
Kim <i>et al.</i> [183]		-	VGG-16	<u>45.36</u>	41.30	55.82
Chen <i>et al.</i> [190]		-	Darknet-53	43.25	<u>46.99</u>	35.84

has examined a wide array of methodologies, including deep learning-based, feature-based, and hybrid approaches, which have demonstrated promising results in improving pedestrian detection performance in challenging lighting scenarios. By delving into the current landscape of low-light pedestrian detection, this work contributes to advancing more secure and dependable autonomous driving systems and other applications related to pedestrian safety. It has also identified ongoing research directions in the field and highlighted potential zones that warrant further research and investigation. The insights provided in this paper aim to inform and inspire future work, ultimately driving innovation and progress in the domain of pedestrian detection under adverse conditions.

REFERENCES

- [1] A. Boukerche and M. Sha, "Design guidelines on deep learning-based pedestrian detection methods for supporting autonomous vehicles," *ACM Computing Surveys (CSUR)*, vol. 54, no. 6, pp. 1–36, 2021.
- [2] M. Sha and A. Boukerche, "Performance evaluation of cnn-based pedestrian detectors for autonomous vehicles," *Ad Hoc Networks*, vol. 128, p. 102784, 2022.
- [3] S. Kim, S. Kwak, and B. C. Ko, "Fast pedestrian detection in surveillance video based on soft target training of shallow random forest," *IEEE Access*, vol. 7, pp. 12415–12426, 2019.
- [4] O. M. Oluyide, J.-R. Tapamo, and T. M. Walingo, "Automatic dynamic range adjustment for pedestrian detection in thermal (infrared) surveillance videos," *Sensors*, vol. 22, no. 5, p. 1728, 2022.
- [5] G. Oltean, L. Ivanciu, and H. Balea, "Pedestrian detection and behaviour characterization for video surveillance systems," in *2019 IEEE 25th International Symposium for Design and Technology in Electronic Packaging (SIITME)*. IEEE, 2019, pp. 256–259.
- [6] M. Zou, J. Yu, B. Lu, W. Chi, and L. Sun, "Active pedestrian detection for excavator robots based on multi-sensor fusion," in *2022 IEEE International Conference on Real-time Computing and Robotics (RCAR)*. IEEE, 2022, pp. 255–260.
- [7] Z. Zhao, X. Qi, Y. Zhao, J. Zhang, W. Wang, and X. Yang, "Pedestrian detection and tracking based on 2d lidar and rgb-d camera," in *Proceedings of the 2022 3rd International Conference on Control, Robotics and Intelligent System*, 2022, pp. 7–14.
- [8] L. Pang, Z. Cao, J. Yu, S. Liang, X. Chen, and W. Zhang, "An efficient 3d pedestrian detector with calibrated rgb camera and 3d lidar," in *2019 IEEE International Conference on Robotics and Biomimetics (ROBIO)*. IEEE, 2019, pp. 2902–2907.
- [9] U. Gawande, K. Hajari, and Y. Golhar, "Pedestrian detection and tracking in video surveillance system: issues, comprehensive review, and challenges," *Recent Trends in Computational Intelligence*, pp. 1–24, 2020.
- [10] W. Wang, X. Chang, J. Yang, and G. Xu, "Lidar-based dense pedestrian detection and tracking," *Applied Sciences*, vol. 12, no. 4, p. 1799, 2022.
- [11] B. Ghari, A. Tourani, and A. Shahbahrani, "A robust pedestrian detection approach for autonomous vehicles," in *Proceedings of the 2022 8th Iranian Conference on Signal Processing and Intelligent Systems*. IEEE, 2022, pp. 1–5.
- [12] Y. Jiang, G. Tong, H. Yin, and N. Xiong, "A pedestrian detection method based on genetic algorithm for optimize xgboost training parameters," *IEEE Access*, vol. 7, pp. 118310–118321, 2019.
- [13] G. Li, Y. Yang, and X. Qu, "Deep learning approaches on pedestrian detection in hazy weather," *IEEE Transactions on Industrial Electronics*, vol. 67, no. 10, pp. 8889–8899, 2019.
- [14] L. Barba-Guaman, J. Eugenio Naranjo, and A. Ortiz, "Deep learning framework for vehicle and pedestrian detection in rural roads on an embedded gpu," *Electronics*, vol. 9, no. 4, p. 589, 2020.
- [15] G. L. Hung, M. S. B. Sahimi, H. Samma, T. A. Almohamad, and B. Lahasan, "Faster r-cnn deep learning model for pedestrian detection from drone images," *SN Computer Science*, vol. 1, pp. 1–9, 2020.
- [16] Z. Ahmed, R. Iniyavan *et al.*, "Enhanced vulnerable pedestrian detection using deep learning," in *2019 International Conference on Communication and Signal Processing (ICCSP)*. IEEE, 2019, pp. 0971–0974.
- [17] L. Chen, S. Lin, X. Lu, D. Cao, H. Wu, C. Guo, C. Liu, and F.-Y. Wang, "Deep neural network based vehicle and pedestrian detection for autonomous driving: A survey," *IEEE Transactions on Intelligent Transportation Systems*, vol. 22, no. 6, pp. 3234–3246, 2021.
- [18] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [19] Y.-L. Hou, Y. Song, X. Hao, Y. Shen, M. Qian, and H. Chen, "Multispectral pedestrian detection based on deep convolutional neural networks," *Infrared Physics & Technology*, vol. 94, pp. 69–77, 2018.
- [20] S. Iftikhar, Z. Zhang, M. Asim, A. Muthanna, A. Koucheryavy, and A. A. Abd El-Latif, "Deep learning-based pedestrian detection in autonomous vehicles: Substantial issues and challenges," *Electronics*, vol. 11, no. 21, p. 3551, 2022.
- [21] W. Chen, Y. Zhu, Z. Tian, F. Zhang, and M. Yao, "Occlusion and multi-scale pedestrian detection a review," *Array*, p. 100318, 2023.
- [22] F. Li, X. Li, Q. Liu, and Z. Li, "Occlusion handling and multi-scale pedestrian detection based on deep learning: a review," *IEEE Access*, vol. 10, pp. 19937–19957, 2022.
- [23] L. Neumann, M. Karg, S. Zhang, C. Scharfenberger, E. Piegert, S. Mistr, O. Prokofyeva, R. Thiel, A. Vedaldi, A. Zisserman *et al.*, "Nightowls: A pedestrians at night dataset," in *Computer Vision—ACCV 2018: 14th Asian Conference on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers, Part I 14*. Springer, 2019, pp. 691–705.
- [24] X. Jia, C. Zhu, M. Li, W. Tang, and W. Zhou, "Llvp: A visible-infrared paired dataset for low-light vision," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 3496–3504.
- [25] S. Hwang, J. Park, N. Kim, Y. Choi, and I. So Kweon, "Multispectral pedestrian detection: Benchmark dataset and baseline," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1037–1045.
- [26] Z. Xu, J. Zhuang, Q. Liu, J. Zhou, and S. Peng, "Benchmarking a large-scale fir dataset for on-road pedestrian detection," *Infrared Physics & Technology*, vol. 96, pp. 199–208, 2019.
- [27] H. Zhang, E. Fromont, S. Lefevre, and B. Avignon, "Multispectral fusion for object detection with cyclic fuse-and-refine blocks," in *2020 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2020, pp. 276–280.
- [28] Y. Socarrás, S. Ramos, D. Vázquez, A. M. López, and T. Gevers, "Adapting pedestrian detection from synthetic to far infrared images," in *ICCV Workshops*, vol. 3, 2013.
- [29] A. González, Z. Fang, Y. Socarrás, J. Serrat, D. Vázquez, J. Xu, and A. M. López, "Pedestrian detection at day/night time with visible and fir cameras: A comparison," *Sensors*, vol. 16, no. 6, p. 820, 2016.
- [30] P. Tumas, A. Nowosielski, and A. Serackis, "Pedestrian detection in severe weather conditions," *IEEE Access*, vol. 8, pp. 62775–62784, 2020.
- [31] J. W. Davis and M. A. Keck, "A two-stage template approach to person detection in thermal imagery," in *2005 Seventh IEEE Workshops on Applications of Computer Vision (WACV/MOTION'05)-Volume 1*, vol. 1. IEEE, 2005, pp. 364–369.
- [32] K. Takumi, K. Watanabe, Q. Ha, A. Tejero-De-Pablos, Y. Ushiku, and T. Harada, "Multispectral object detection for autonomous vehicles," in *Proceedings of the on Thematic Workshops of ACM Multimedia 2017*, 2017, pp. 35–43.
- [33] M. Jeong, B. C. Ko, and J.-Y. Nam, "Early detection of sudden pedestrian crossing for safe driving during summer nights," *IEEE transactions on circuits and systems for video technology*, vol. 27, no. 6, pp. 1368–1380, 2016.
- [34] D. Olmeda, C. Premebida, U. Nunes, J. M. Armingol, and A. de la Escalera, "Pedestrian detection in far infrared images," *Integrated Computer-Aided Engineering*, vol. 20, no. 4, pp. 347–360, 2013.
- [35] Z. Wu, N. Fuller, D. Thériault, and M. Betke, "A thermal infrared video benchmark for visual analysis," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2014, pp. 201–208.
- [36] E. Gebhardt and M. Wolf, "Camel dataset for visual and thermal infrared multiple object detection and tracking," in *2018 15th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*. IEEE, 2018, pp. 1–6.
- [37] M. A. Farooq, W. Shariff, and P. Corcoran, "Evaluation of thermal imaging on embedded gpu platforms for application in vehicular assistance systems," *arXiv preprint arXiv:2201.01661*, 2022.
- [38] P. Karol, P. Paweł, and D. Adam, "Video processing algorithms for detection of pedestrians," *CMST*, vol. 21, no. 3, pp. 141–150, 2015.

- [39] A. Torabi, G. Massé, and G.-A. Bilodeau, "An iterative integrated framework for thermal-visible image registration, sensor fusion, and people tracking for video surveillance applications," *Computer Vision and Image Understanding*, vol. 116, no. 2, pp. 210–221, 2012.
- [40] T. Kim and S. Kim, "Pedestrian detection at night time in fir domain: Comprehensive study about temperature and brightness and new benchmark," *Pattern Recognition*, vol. 79, pp. 44–54, 2018.
- [41] "Flir thermal dataset for algorithm training," <https://www.flir.in/oem/adas/adas-dataset-form/>.
- [42] A. Nowosielski, K. Malecki, P. Forczmański, A. Smoliński, and K. Krzywicki, "Embedded night-vision system for pedestrian detection," *IEEE Sensors Journal*, vol. 20, no. 16, pp. 9293–9304, 2020.
- [43] J. Kim, "Pedestrian detection and distance estimation using thermal camera in night time," in *2019 International Conference on Artificial Intelligence in Information and Communication (ICAIC)*. IEEE, 2019, pp. 463–466.
- [44] J. B. Kim, "Detection of direction indicators on road surfaces using inverse perspective mapping and nn," *J. Inf. Process. Korean*, vol. 4, pp. 201–208, 2015.
- [45] D. Zhou, S. Qiu, Y. Song, and K. Xia, "A pedestrian extraction algorithm based on single infrared image," *Infrared Physics & Technology*, vol. 105, p. 103236, 2020.
- [46] A. B. Khalifa, I. Alouani, M. A. Mahjoub, and N. E. B. Amara, "Pedestrian detection using a moving camera: A novel framework for foreground detection," *Cognitive Systems Research*, vol. 60, pp. 77–96, 2020.
- [47] A. R. Shahzad and A. Jalal, "A smart surveillance system for pedestrian tracking and counting using template matching," in *2021 International Conference on Robotics and Automation in Industry (ICRAI)*. IEEE, 2021, pp. 1–6.
- [48] Y. Cai, Z. Liu, H. Wang, and X. Sun, "Saliency-based pedestrian detection in far infrared images," *IEEE Access*, vol. 5, pp. 5013–5019, 2017.
- [49] H. Zhang, H. Xu, X. Tian, J. Jiang, and J. Ma, "Image fusion meets deep learning: A survey and perspective," *Information Fusion*, vol. 76, pp. 323–336, 2021.
- [50] J. Nataprawira, Y. Gu, I. Goncharenko, and S. Kamijo, "Pedestrian detection on multispectral images in different lighting conditions," in *2021 IEEE International Conference on Consumer Electronics (ICCE)*. IEEE, 2021, pp. 1–5.
- [51] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018.
- [52] J. Nataprawira, Y. Gu, I. Goncharenko, and S. Kamijo, "Pedestrian detection using multispectral images and a deep neural network," *Sensors*, vol. 21, no. 7, p. 2536, 2021.
- [53] Y. Yang, K. Xu, and K. Wang, "Cascaded information enhancement and cross-modal attention feature fusion for multispectral pedestrian detection," *Frontiers in Physics*, vol. 11, p. 1121311, 2023.
- [54] Y. Zhang, Z. Yin, L. Nie, and S. Huang, "Attention based multi-layer fusion of multispectral images for pedestrian detection," *IEEE Access*, vol. 8, pp. 165 071–165 084, 2020.
- [55] Q. Jiang, J. Dai, T. Rui, F. Shao, J. Wang, and G. Lu, "Attention-based cross-modality feature complementation for multispectral pedestrian detection," *IEEE Access*, vol. 10, pp. 53 797–53 809, 2022.
- [56] L. Fu, W.-b. Gu, Y.-b. Ai, W. Li, and D. Wang, "Adaptive spatial pixel-level feature fusion network for multispectral pedestrian detection," *Infrared Physics & Technology*, vol. 116, p. 103770, 2021.
- [57] Q. Deng, W. Tian, Y. Huang, L. Xiong, and X. Bi, "Pedestrian detection by fusion of rgb and infrared images in low-light environment," in *2021 IEEE 24th International Conference on Information Fusion (FUSION)*. IEEE, 2021, pp. 1–8.
- [58] J. Liu, S. Zhang, S. Wang, and D. N. Metaxas, "Multispectral deep neural networks for pedestrian detection," *arXiv preprint arXiv:1611.02644*, 2016.
- [59] Y. Cao, J. Bin, J. Hamari, E. Blasch, and Z. Liu, "Multimodal object detection by channel switching and spatial attention," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 403–411.
- [60] X. Yang, Y. Qian, H. Zhu, C. Wang, and M. Yang, "Baanet: Learning bi-directional adaptive attention gates for multispectral pedestrian detection," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 2920–2926.
- [61] H. Zhang, E. Fromont, S. Lefèvre, and B. Avignon, "Guided attentive feature fusion for multispectral pedestrian detection," in *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, 2021, pp. 72–80.
- [62] F. Qingyun, H. Dapeng, and W. Zhaokui, "Cross-modality fusion transformer for multispectral object detection," *arXiv preprint arXiv:2111.00273*, 2021.
- [63] D. Pei, M. Jing, H. Liu, F. Sun, and L. Jiang, "A fast retinanet fusion framework for multi-spectral pedestrian detection," *Infrared Physics & Technology*, vol. 105, p. 103178, 2020.
- [64] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.
- [65] L. Ding, Y. Wang, R. Laganiere, D. Huang, and S. Fu, "Convolutional neural networks for multispectral pedestrian detection," *Signal Processing: Image Communication*, vol. 82, p. 115764, 2020.
- [66] J. Dai, Y. Li, K. He, and J. Sun, "R-fcn: Object detection via region-based fully convolutional networks," *Advances in neural information processing systems*, vol. 29, 2016.
- [67] J.-S. Yun, S.-H. Park, and S. B. Yoo, "Infusion-net: Inter-and intra-weighted cross-fusion network for multispectral object detection," *Mathematics*, vol. 10, no. 21, p. 3966, 2022.
- [68] C. Bao, J. Cao, Q. Hao, Y. Cheng, Y. Ning, and T. Zhao, "Dual-yolo architecture from infrared and visible images for object detection," *Sensors*, vol. 23, no. 6, p. 2934, 2023.
- [69] C.-Y. Wang, A. Bochkovskiy, and H.-Y. M. Liao, "Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 7464–7475.
- [70] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points. arxiv 2019," *arXiv preprint arXiv:1904.07850*, vol. 448, 1904.
- [71] X. Zuo, Z. Wang, J. Shen, and W. Yang, "Improving multispectral pedestrian detection with scale-aware permutation attention and adjacent feature aggregation," *IET Computer Vision*, 2022.
- [72] Y. Cao, D. Guan, Y. Wu, J. Yang, Y. Cao, and M. Y. Yang, "Box-level segmentation supervised deep neural networks for accurate and real-time multispectral pedestrian detection," *ISPRS journal of photogrammetry and remote sensing*, vol. 150, pp. 70–79, 2019.
- [73] L. Tang, J. Yuan, H. Zhang, X. Jiang, and J. Ma, "Piafusion: A progressive infrared and visible image fusion network based on illumination aware," *Information Fusion*, vol. 83, pp. 79–92, 2022.
- [74] K. Roszyk, M. R. Nowicki, and P. Skrzypczyński, "Adopting the yolov4 architecture for low-latency multispectral pedestrian detection in autonomous driving," *Sensors*, vol. 22, no. 3, p. 1082, 2022.
- [75] P. Peng, T. Xu, B. Huang, and J. Li, "Hafnet: Hierarchical attentive fusion network for multispectral pedestrian detection," *Remote Sensing*, vol. 15, no. 8, p. 2041, 2023.
- [76] R. Yadav, A. Samir, H. Rashed, S. Yogamani, and R. Dahyot, "Cnn based color and thermal image fusion for object detection in automated driving," *Irish Machine Vision and Image Processing*, 2020.
- [77] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *Advances in neural information processing systems*, vol. 28, 2015.
- [78] L. Zhang, Z. Liu, S. Zhang, X. Yang, H. Qiao, K. Huang, and A. Husain, "Cross-modality interactive attention network for multispectral pedestrian detection," *Information Fusion*, vol. 50, pp. 20–29, 2019.
- [79] R. Hu, T. Rui, Y. Ouyang, J. Wang, Q. Jiang, and Y. Du, "Dmffnet: Dual-mode multi-scale feature fusion-based pedestrian detection method," *IEEE Access*, 2022.
- [80] A. Howard, M. Sandler, G. Chu, L.-C. Chen, B. Chen, M. Tan, W. Wang, Y. Zhu, R. Pang, V. Vasudevan *et al.*, "Searching for mobilenetv3," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 1314–1324.
- [81] Z. Cao, H. Yang, J. Zhao, S. Guo, and L. Li, "Attention fusion for one-stage multispectral pedestrian detection," *Sensors*, vol. 21, no. 12, p. 4184, 2021.
- [82] J. Kim, J. Choi, Y. Kim, J. Koh, C. C. Chung, and J. W. Choi, "Robust camera lidar sensor fusion via deep gated information fusion network," in *2018 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2018, pp. 1620–1625.
- [83] Y. Zheng, I. H. Izzat, and S. Ziaee, "Gfd-ssd: gated fusion double ssd for multispectral pedestrian detection," *arXiv preprint arXiv:1903.06999*, 2019.
- [84] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*. Springer, 2016, pp. 21–37.
- [85] Q. Wang, Y. Chi, T. Shen, J. Song, Z. Zhang, and Y. Zhu, "Improving rgb-infrared object detection by reducing cross-modality redundancy," *Remote Sensing*, vol. 14, no. 9, p. 2020, 2022.

- [86] Y. Xing, S. Wang, G. Liang, Q. Li, X. Zhang, S. Zhang, and Y. Zhang, "Multispectral pedestrian detection via reference box constrained cross attention and modality balanced optimization," *arXiv preprint arXiv:2302.00290*, 2023.
- [87] B. Khalid, A. M. Khan, M. U. Akram, and S. Batool, "Person detection by fusion of visible and thermal images using convolutional neural network," in *2019 2nd International Conference on Communication, Computing and Digital systems (C-CODE)*. IEEE, 2019, pp. 143–148.
- [88] B. Montenegro and M. Flores-Calero, "Pedestrian detection at daytime and nighttime conditions based on yolo-v5," *Ingenius. Revista de Ciencia y Tecnología*, no. 27, pp. 85–95, 2022.
- [89] X. Song, S. Gao, and C. Chen, "A multispectral feature fusion network for robust pedestrian detection," *Alexandria Engineering Journal*, vol. 60, no. 1, pp. 73–85, 2021.
- [90] L. Ding, Y. Wang, R. Laganieri, D. Huang, X. Luo, and H. Zhang, "A robust and fast multispectral pedestrian detection deep network," *Knowledge-Based Systems*, vol. 227, p. 106990, 2021.
- [91] K. Park, S. Kim, and K. Sohn, "Unified multi-spectral pedestrian detection based on probabilistic fusion networks," *Pattern Recognition*, vol. 80, pp. 143–155, 2018.
- [92] A. Bochkovski, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," *arXiv preprint arXiv:2004.10934*, 2020.
- [93] H.-T. Chan, P.-T. Tsai, and C.-H. Hsia, "Multispectral pedestrian detection via two-stream yolo with complementarity fusion for autonomous driving," in *2023 IEEE 3rd International Conference on Electronic Communications, Internet of Things and Big Data (ICEIB)*. IEEE, 2023, pp. 313–316.
- [94] Y.-T. Chen, J. Shi, Z. Ye, C. Mertz, D. Ramanan, and S. Kong, "Multimodal object detection via probabilistic ensembling," in *European Conference on Computer Vision*. Springer, 2022, pp. 139–158.
- [95] Z. A. Shaikh, D. Van Hamme, P. Veelaert, and W. Philips, "Probabilistic fusion for pedestrian detection from thermal and colour images," *Sensors*, vol. 22, no. 22, p. 8637, 2022.
- [96] Y. Zhuang, Z. Pu, J. Hu, and Y. Wang, "Illumination and temperature-aware multispectral networks for edge-computing-enabled pedestrian detection," *IEEE Transactions on Network Science and Engineering*, vol. 9, no. 3, pp. 1282–1295, 2021.
- [97] C.-H. Hsia, H.-C. Peng, and H.-T. Chan, "All-weather pedestrian detection based on double-stream multispectral network," *Electronics*, vol. 12, no. 10, p. 2312, 2023.
- [98] C. Li, D. Song, R. Tong, and M. Tang, "Illumination-aware faster r-cnn for robust multispectral pedestrian detection," *Pattern Recognition*, vol. 85, pp. 161–171, 2019.
- [99] G. Li, W. Lai, and X. Qu, "Pedestrian detection based on light perception fusion of visible and thermal images," *Optics & Laser Technology*, vol. 156, p. 108466, 2022.
- [100] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.
- [101] Y. Cao, X. Luo, J. Yang, Y. Cao, and M. Y. Yang, "Locality guided cross-modal feature aggregation and pixel-level fusion for multispectral pedestrian detection," *Information Fusion*, vol. 88, pp. 1–11, 2022.
- [102] A. Wolpert, M. Teutsch, M. S. Sarfraz, and R. Stiefelwagen, "Anchor-free small-scale multispectral pedestrian detection," *arXiv preprint arXiv:2008.08418*, 2020.
- [103] J. Hu, Y. Zhao, and X. Zhang, "Application of transfer learning in infrared pedestrian detection," in *2020 IEEE 5th International Conference on Image, Vision and Computing (ICIVC)*. IEEE, 2020, pp. 1–4.
- [104] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.
- [105] M. Vandersteegen, K. Van Beeck, and T. Goedemé, "Real-time multispectral pedestrian detection with a single-pass deep neural network," in *Image Analysis and Recognition: 15th International Conference, ICIAR 2018, Póvoa de Varzim, Portugal, June 27–29, 2018, Proceedings 15*. Springer, 2018, pp. 419–426.
- [106] J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7263–7271.
- [107] S. Geng, "Infrared image pedestrian target detection based on yolov3 and migration learning," *arXiv preprint arXiv:2012.11185*, 2020.
- [108] Z. Zheng, P. Wang, W. Liu, J. Li, R. Ye, and D. Ren, "Distance-iou loss: Faster and better learning for bounding box regression," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 07, 2020, pp. 12 993–13 000.
- [109] T. Guo, C. P. Huynh, and M. Solh, "Domain-adaptive pedestrian detection in thermal images," in *2019 IEEE international conference on image processing (ICIP)*. IEEE, 2019, pp. 1660–1664.
- [110] M. Kieu, A. D. Bagdanov, M. Bertini, and A. Del Bimbo, "Task-conditioned domain adaptation for pedestrian detection in thermal imagery," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXII 16*. Springer, 2020, pp. 546–562.
- [111] —, "Domain adaptation for privacy-preserving pedestrian detection in thermal imagery," in *Image Analysis and Processing—ICIAP 2019: 20th International Conference, Trento, Italy, September 9–13, 2019, Proceedings, Part II 20*. Springer, 2019, pp. 203–213.
- [112] M. Kieu, A. D. Bagdanov, and M. Bertini, "Bottom-up and layerwise domain adaptation for pedestrian detection in thermal images," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 17, no. 1, pp. 1–19, 2021.
- [113] M. Krišto, M. Ivacic-Kos, and M. Pobar, "Thermal object detection in difficult weather conditions using yolo," *IEEE access*, vol. 8, pp. 125 459–125 476, 2020.
- [114] K. Fritz, D. König, U. Klauck, and M. Teutsch, "Generalization ability of region proposal networks for multispectral person detection," in *Automatic Target Recognition XXIX*, vol. 10988. SPIE, 2019, pp. 222–235.
- [115] P. Dollar, C. Wojek, B. Schiele, and P. Perona, "Pedestrian detection: An evaluation of the state of the art," *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 4, pp. 743–761, 2011.
- [116] S. Zhang, R. Benenson, and B. Schiele, "Citypersons: A diverse dataset for pedestrian detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3213–3221.
- [117] Q. Ha, K. Watanabe, T. Karasawa, Y. Ushiku, and T. Harada, "Mfnet: Towards real-time semantic segmentation for autonomous vehicles with multi-spectral scenes," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 5108–5115.
- [118] V. Vs, D. Poster, S. You, S. Hu, and V. M. Patel, "Meta-uda: Unsupervised domain adaptive thermal object detection using meta-learning," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 1412–1423.
- [119] C. Lyu, P. Heyer, A. Munir, L. Platasa, C. Micheloni, B. Goossens, and W. Philips, "Visible-thermal pedestrian detection via unsupervised transfer learning," in *2021 the 5th International Conference on Innovation in Artificial Intelligence*, 2021, pp. 158–163.
- [120] F. Munir, S. Azam, and M. Jeon, "Sstm: Self-supervised domain adaptation thermal object detection for autonomous driving," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 206–213.
- [121] C. Lyu, P. Heyer, B. Goossens, and W. Philips, "An unsupervised transfer learning framework for visible-thermal pedestrian detection," *Sensors*, vol. 22, no. 12, p. 4416, 2022.
- [122] Y. Cao, D. Guan, W. Huang, J. Yang, Y. Cao, and Y. Qiao, "Pedestrian detection with unsupervised multispectral feature learning using deep neural networks," *information fusion*, vol. 46, pp. 206–217, 2019.
- [123] D. Guan, X. Luo, Y. Cao, J. Yang, Y. Cao, G. Vosselman, and M. Ying Yang, "Unsupervised domain adaptation for multispectral pedestrian detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.
- [124] M. A. Marnissi, H. Fradi, A. Sahbani, and N. E. B. Amara, "Unsupervised thermal-to-visible domain adaptation method for pedestrian detection," *Pattern Recognition Letters*, vol. 153, pp. 222–231, 2022.
- [125] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.
- [126] M. Hniewa, A. Rahimpour, J. Miller, D. Upadhyay, and H. Radha, "Cross modality knowledge distillation for robust pedestrian detection in low light and adverse weather conditions," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.
- [127] H. Zhang, E. Fromont, S. Lefèvre, and B. Avignon, "Low-cost multispectral scene analysis with modality distillation," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 803–812.
- [128] T. Liu, K.-M. Lam, R. Zhao, and G. Qiu, "Deep cross-modal representation learning and distillation for illumination-invariant pedestrian detection," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 1, pp. 315–329, 2021.

- [129] J. U. Kim, S. Park, and Y. M. Ro, "Robust small-scale pedestrian detection with cued recall via memory learning," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 3050–3059.
- [130] —, "Towards versatile pedestrian detector with multisensory-matching and multispectral recalling memory," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 1, 2022, pp. 1157–1165.
- [131] M. A. Marnissi, H. Fradi, A. Sahbani, and N. E. B. Amara, "Thermal image enhancement using generative adversarial network for pedestrian detection," in *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 2021, pp. 6509–6516.
- [132] Y. Sun, Y. Shao, G. Yang, and H. Xie, "A method of infrared image pedestrian detection with improved yolov3 algorithm," *American Journal of Optics and Photonics*, vol. 9, no. 3, pp. 32–38, 2021.
- [133] J. Yu, Y. Fan, J. Yang, N. Xu, Z. Wang, X. Wang, and T. Huang, "Wide activation for efficient and accurate image super-resolution," *arXiv preprint arXiv:1808.08718*, 2018.
- [134] M. A. Marnissi and A. Fathallah, "Gan-based vision transformer for high-quality thermal image enhancement," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 817–825.
- [135] L. Tang, X. Xiang, H. Zhang, M. Gong, and J. Ma, "Divfusion: Darkness-free infrared and visible image fusion," *Information Fusion*, vol. 91, pp. 477–493, 2023.
- [136] G. Li, S. Zhang, and J. Yang, "Nighttime pedestrian detection based on feature attention and transformation," in *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 2021, pp. 9180–9187.
- [137] C. Cui, J. Xie, and Y. Yang, "Bright channel prior attention for multispectral pedestrian detection," *arXiv preprint arXiv:2305.12845*, 2023.
- [138] Y. Chen and H. Shin, "Pedestrian detection at night in infrared images using an attention-guided encoder-decoder convolutional neural network," *Applied Sciences*, vol. 10, no. 3, p. 809, 2020.
- [139] H. Patel, K. Prajapati, A. Sarvaiya, K. Upla, K. Raja, R. Ramachandra, and C. Busch, "Depthwise convolution for compact object detector in nighttime images," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 379–389.
- [140] S. Li, Y. Li, Y. Li, M. Li, and X. Xu, "Yolo-firi: Improved yolov5 for infrared image object detection," *IEEE access*, vol. 9, pp. 141 861–141 875, 2021.
- [141] H. Li and X.-J. Wu, "Densefuse: A fusion approach to infrared and visible images," *IEEE Transactions on Image Processing*, vol. 28, no. 5, pp. 2614–2623, 2018.
- [142] Y. Pang, J. Lin, T. Qin, and Z. Chen, "Image-to-image translation: Methods and applications," *IEEE Transactions on Multimedia*, vol. 24, pp. 3859–3881, 2021.
- [143] F. Luo, Y. Li, G. Zeng, P. Peng, G. Wang, and Y. Li, "Thermal infrared image colorization for nighttime driving scenes with top-down guided attention," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 9, pp. 15 808–15 823, 2022.
- [144] A. Dangle, R. Mundada, S. Gore, J. Shrugare, and H. Dalal, "Enhanced colorization of thermal images for pedestrian detection using deep convolutional neural networks," *Procedia Computer Science*, vol. 218, pp. 2091–2101, 2023.
- [145] S. Yang, M. Sun, X. Lou, H. Yang, and H. Zhou, "An unpaired thermal infrared image translation method using gma-cycleGAN," *Remote Sensing*, vol. 15, no. 3, p. 663, 2023.
- [146] C. Devaguptapu, N. Akolekar, M. M. Sharma, and V. N. Balasubramanian, "Borrow from anywhere: Pseudo multi-modal object detection in thermal imagery," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.
- [147] M. Kieu, L. Berlincioni, L. Galteri, M. Bertini, A. D. Bagdanov, and A. Del Bimbo, "Robust pedestrian detection in thermal imagery using synthesized images," in *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 2021, pp. 8804–8811.
- [148] X. Mao, Q. Li, H. Xie, R. Y. Lau, and Z. Wang, "Multi-class generative adversarial networks with the l2 loss function," *arXiv preprint arXiv:1611.04076*, vol. 5, pp. 1057–1149, 2016.
- [149] F. Altay and S. Velipasalar, "The use of thermal cameras for pedestrian detection," *IEEE Sensors Journal*, vol. 22, no. 12, pp. 11 489–11 498, 2022.
- [150] D. Ghose, S. M. Desai, S. Bhattacharya, D. Chakraborty, M. Fiterau, and T. Rahman, "Pedestrian detection in thermal images using saliency maps," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.
- [151] N. Liu, J. Han, and M.-H. Yang, "Picanet: Learning pixel-wise contextual attention for saliency detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3089–3098.
- [152] Z. Deng, X. Hu, L. Zhu, X. Xu, J. Qin, G. Han, and P.-A. Heng, "R3net: Recurrent residual refinement network for saliency detection," in *Proceedings of the 27th International Joint Conference on Artificial Intelligence*. AAAI Press Menlo Park, CA, USA, 2018, pp. 684–690.
- [153] M. A. Marnissi, I. Hattab, H. Fradi, A. Sahbani, and N. E. B. Amara, "Bispectral pedestrian detection augmented with saliency maps using transformer," in *VISIGRAPP (5: VISAPP)*, 2022, pp. 275–284.
- [154] Y. Zhao, J. Cheng, W. Zhou, C. Zhang, and X. Pan, "Infrared pedestrian detection with converted temperature map," in *2019 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*. IEEE, 2019, pp. 2025–2031.
- [155] K. Oksuz, B. C. Cam, S. Kalkan, and E. Akbas, "Imbalance problems in object detection: A review," *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 10, pp. 3388–3415, 2020.
- [156] K. Zhou, L. Chen, and X. Cao, "Improving multispectral pedestrian detection by addressing modality imbalance problems," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVIII 16*. Springer, 2020, pp. 787–803.
- [157] A. Das, S. Das, G. Sistu, J. Horgan, U. Bhattacharya, E. Jones, M. Glavin, and C. Eising, "Revisiting modality imbalance in multi-modal pedestrian detection," *arXiv preprint arXiv:2302.12589*, 2023.
- [158] L. Gross, "Logarithmic sobolev inequalities," *American Journal of Mathematics*, vol. 97, no. 4, pp. 1061–1083, 1975.
- [159] W. Li, "Infrared image pedestrian detection via yolo-v3," in *2021 IEEE 5th Advanced Information Technology, Electronic and Automation Control Conference (IAEAC)*, vol. 5. IEEE, 2021, pp. 1052–1055.
- [160] K. Dasgupta, A. Das, S. Das, U. Bhattacharya, and S. Yogamani, "Spatio-contextual deep network-based multimodal pedestrian detection for autonomous driving," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 9, pp. 15 940–15 950, 2022.
- [161] P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Lio, Y. Bengio et al., "Graph attention networks," *stat*, vol. 1050, no. 20, pp. 10–48 550, 2017.
- [162] S. Bell, C. L. Zitnick, K. Bala, and R. Girshick, "Inside-outside net: Detecting objects in context with skip pooling and recurrent neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2874–2883.
- [163] L. Zhang, X. Zhu, X. Chen, X. Yang, Z. Lei, and Z. Liu, "Weakly aligned cross-modal learning for multispectral pedestrian detection," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 5127–5137.
- [164] M. Kim, S. Joung, K. Park, S. Kim, and K. Sohn, "Unpaired cross-spectral pedestrian detection via adversarial feature learning," in *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 1650–1654.
- [165] J. U. Kim, S. Park, and Y. M. Ro, "Uncertainty-guided cross-modal learning for robust multispectral pedestrian detection," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 3, pp. 1510–1523, 2021.
- [166] N. Wanchaitanawong, M. Tanaka, T. Shibata, and M. Okutomi, "Multi-modal pedestrian detection with large misalignment based on modal-wise regression and multi-modal iou," in *2021 17th International Conference on Machine Vision and Applications (MVA)*. IEEE, 2021, pp. 1–6.
- [167] M. Crawshaw, "Multi-task learning with deep neural networks: A survey," *arXiv preprint arXiv:2009.09796*, 2020.
- [168] Y. Wang, T. Lu, Y. Zhang, W. Fang, Y. Wu, and Z. Wang, "Cross-task feature alignment for seeing pedestrians in the dark," *Neurocomputing*, vol. 462, pp. 282–293, 2021.
- [169] D. Guan, Y. Cao, J. Yang, Y. Cao, and M. Y. Yang, "Fusion of multispectral data through illumination-aware deep neural networks for pedestrian detection," *Information Fusion*, vol. 50, pp. 148–157, 2019.
- [170] X. Dai, J. Hu, H. Zhang, A. Shitu, C. Luo, A. Osman, S. Sfarra, and Y. Duan, "Multi-task faster r-cnn for nighttime pedestrian detection and distance estimation," *Infrared Physics & Technology*, vol. 115, p. 103694, 2021.
- [171] Z. Cao, H. Yang, J. Zhao, X. Pan, L. Zhang, and Z. Liu, "A new region proposal network for far-infrared pedestrian detection," *IEEE Access*, vol. 7, pp. 135 023–135 030, 2019.
- [172] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the*

- IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125.
- [173] C. Li, D. Song, R. Tong, and M. Tang, “Multispectral pedestrian detection via simultaneous detection and segmentation,” *arXiv preprint arXiv:1808.04818*, 2018.
 - [174] Y.-Y. Chen, S.-Y. Jhong, G.-Y. Li, and P.-H. Chen, “Thermal-based pedestrian detection using faster r-cnn and region decomposition branch,” in *2019 International Symposium on Intelligent Signal Processing and Communication Systems (ISPACS)*. IEEE, 2019, pp. 1–2.
 - [175] W. Liu, S. Liao, W. Ren, W. Hu, and Y. Yu, “High-level semantic feature detection: A new perspective for pedestrian detection,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 5187–5196.
 - [176] Y. Song, M. Li, X. Qiu, W. Du, and J. Feng, “Full-time infrared feature pedestrian detection based on csp network,” in *2020 International Conference on Intelligent Transportation, Big Data & Smart City (ICITBS)*. IEEE, 2020, pp. 516–518.
 - [177] Z. Xu, C.-M. Wong, C.-C. Wong, and Q. Liu, “Ground plane context aggregation network for day-and-night on vehicular pedestrian detection,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 22, no. 10, pp. 6395–6406, 2020.
 - [178] X. Dai, Y. Duan, J. Hu, S. Liu, C. Hu, Y. He, D. Chen, C. Luo, and J. Meng, “Near infrared nighttime road pedestrians recognition based on convolutional neural network,” *Infrared Physics & Technology*, vol. 97, pp. 25–32, 2019.
 - [179] C. M. Bishop and N. M. Nasrabadi, *Pattern recognition and machine learning*. Springer, 2006, vol. 4, no. 4.
 - [180] M. A. Galarza-Bravo and M. J. Flores-Calero, “Pedestrian detection at night based on faster r-cnn and far infrared images,” in *Intelligent Robotics and Applications: 11th International Conference, ICIRA 2018, Newcastle, NSW, Australia, August 9–11, 2018, Proceedings, Part II*. Springer, 2018, pp. 335–345.
 - [181] R. Kalita, A. K. Talukdar, and K. K. Sarma, “Real-time human detection with thermal camera feed using yolov3,” in *2020 IEEE 17th India Council International Conference (INDICON)*. IEEE, 2020, pp. 1–5.
 - [182] K. N. R. Chebrolu and P. Kumar, “Deep learning based pedestrian detection at all light conditions,” in *2019 International Conference on Communication and Signal Processing (ICCSP)*. IEEE, 2019, pp. 0838–0842.
 - [183] J. H. Kim, G. Batchuluun, and K. R. Park, “Pedestrian detection based on faster r-cnn in nighttime by fusing deep convolutional features of successive images,” *Expert Systems with Applications*, vol. 114, pp. 15–33, 2018.
 - [184] D. König, M. Adam, C. Jarvers, G. Layher, H. Neumann, and M. Teutsch, “Fully convolutional region proposal networks for multispectral person detection,” in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2017, pp. 49–56.
 - [185] P. Tumas, A. Jonkus, and A. Serackis, “Acceleration of hog based pedestrian detection in fir camera video stream,” in *2018 Open Conference of Electrical, Electronic and Information Sciences (eStream)*. IEEE, 2018, pp. 1–4.
 - [186] A. Narayanan, R. D. Kumar, R. RoselinKiruba, and T. S. Sharmila, “Study and analysis of pedestrian detection in thermal images using yolo and svm,” in *2021 Sixth International Conference on Wireless Communications, Signal Processing and Networking (WiSPNET)*. IEEE, 2021, pp. 431–434.
 - [187] D. Xu, W. Ouyang, E. Ricci, X. Wang, and N. Sebe, “Learning cross-modal deep representations for robust pedestrian detection,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5363–5371.
 - [188] P. Dollár, R. Appel, S. Belongie, and P. Perona, “Fast feature pyramids for object detection,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 8, pp. 1532–1545, 2014.
 - [189] H. Choi, S. Kim, K. Park, and K. Sohn, “Multi-spectral pedestrian detection based on accumulated object proposal with fully convolutional networks,” in *2016 23rd International conference on pattern recognition (ICPR)*. IEEE, 2016, pp. 621–626.
 - [190] X. Chen, L. Liu, and X. Tan, “Robust pedestrian detection based on multi-spectral image fusion and convolutional neural networks,” *Electronics*, vol. 11, no. 1, p. 1, 2021.
 - [191] Y. Ma, J. Chen, C. Chen, F. Fan, and J. Ma, “Infrared and visible image fusion using total variation model,” *Neurocomputing*, vol. 202, pp. 12–19, 2016.