

Pedestrian Detection in Thermal Images using Saliency Maps

Debasmita Ghose* Shasvat M. Desai* Sneha Bhattacharya* Deep Chakraborty*
Madalina Fiterau Tauhidur Rahman

College of Information and Computer Sciences, University of Massachusetts, Amherst, MA 01002

{dghose, shasvatmukes, snehabhattac, dchakraborty, mfilterau, trahman}@cs.umass.edu

Abstract

Thermal images are mainly used to detect the presence of people at night or in bad lighting conditions, but perform poorly at daytime. To solve this problem, most state-of-the-art techniques employ a fusion network that uses features from paired thermal and color images. Instead, we propose to augment thermal images with their saliency maps, to serve as an attention mechanism for the pedestrian detector especially during daytime. We investigate how such an approach results in improved performance for pedestrian detection using only thermal images, eliminating the need for paired color images. For our experiments, we train the Faster R-CNN for pedestrian detection and report the added effect of saliency maps generated using static and deep methods (PiCA-Net and R^3 -Net). Our best performing model results in an absolute reduction of miss rate by 13.4% and 19.4% over the baseline in day and night images respectively. We also annotate and release pixel level masks of pedestrians on a subset of the KAIST Multispectral Pedestrian Detection dataset, which is a first publicly available dataset for salient pedestrian detection.

1. Introduction

Detecting the presence and location of pedestrians in a scene is a crucial task for several applications such as video surveillance systems [38] and autonomous driving [13]. Despite the challenges associated with it, such as low resolution and occlusion, pedestrian detection has already been successfully studied widely in color images and videos using state-of-the-art deep learning techniques for object detection and/or semantic segmentation [4, 24, 11, 3]. Color images of reasonable quality are good for detecting pedestrians during the day. Thermal images, however, are very useful in detecting pedestrians in conditions where color images fail, such as nighttime or under bad lighting conditions. This is because at nighttime, thermal cameras cap-

ture humans distinctly as they are warmer than their surrounding objects. During the day however, there are other objects in the surroundings which are as warm as or warmer than humans, making them less distinguishable. Therefore, there appears to be a clear complementary potential between color and thermal images. In order to exploit this complementary potential, there has been a lot of work on building fusion architectures combining color and thermal images [37, 41, 26, 23]. But color-thermal image pairs might not always be available, as they are expensive to collect and need image registration to be completely accurate. Misaligned imagery can also reduce the performance of a detector that leverages multiple data modalities. This motivates us to use only thermal images for the task of pedestrian detection.

To address the challenge of pedestrian detection in thermal images, especially during daytime, we propose the use of saliency maps. Koch and Ullman [21] define saliency at a given location by how different this location is from its surroundings in color, orientation, motion, and depth. Looking for salient objects in a scene can be interpreted as being a visual attention mechanism which illuminates pixels belonging to salient objects in a given scene. We therefore hypothesize that using saliency maps along with thermal images would help us improve the performance of state-of-the-art pedestrian detectors, especially on thermal images captured during the day. To test our hypothesis, we first establish a baseline by training a state-of-the-art object detector (Faster R-CNN [32]) to detect pedestrians solely from thermal images in the KAIST Multispectral Pedestrian dataset [18]. We then train pedestrian detectors on thermal images augmented with their saliency maps generated using static and deep learning techniques (PiCA-Net[28] and R^3 -Net[7]). Our experiments show that the pedestrian detector trained using this augmentation technique outperforms the baseline by a significant margin. Moreover, since deep saliency networks require pixel level annotations of salient objects, we annotate a subset of the KAIST multispectral pedestrian dataset [18] with pixel level masks for pedestrian instances to facilitate research on salient pedestrian detection.

* Authors contributed equally

The key contributions of this paper are as follows:

1. To the best of our knowledge, this is the first paper to show the impact of saliency maps in improving the performance of pedestrian detection in thermal images.
2. We release the first pixel level annotations for a multispectral pedestrian detection dataset and provide saliency detection benchmarks on it using state-of-the-art networks.

The rest of the paper is organized as follows: Section 2 reviews existing work on pedestrian detection in color and multispectral images and methods for saliency detection in images. Section 3 outlines the baseline method for pedestrian detection and our efforts to improve it using saliency maps. We also present a new salient pedestrian detection dataset that we annotated for this purpose. In Section 4 we report implementation details, benchmarks for our novel dataset and evaluate the performance of different techniques qualitatively and quantitatively. Finally, we present our conclusions and future work in Section 5.

2. Related Work

Pedestrian detection. Traditionally, pedestrian detectors involved the use of hand crafted features and algorithms such as ICF [10], ACF [9] and LDCF [31]. Deep learning approaches have however been more successful recently. Zhang *et al.* [44] investigate the performance of the Faster R-CNN [32] for the task of pedestrian detection. Sermanet *et al.* [34] introduce the use of multistage unsupervised features and skip connections for pedestrian detection. Li *et al.* [24] introduce Scale Aware Fast R-CNN which uses built-in sub-networks to detect pedestrians at different scales. In [3], Brazil *et al.* introduce SDS R-CNN which uses joint supervision on pedestrian detection and semantic segmentation to illuminate pedestrians in the frame. This motivates us to use saliency maps as a stronger attention mechanism to illuminate pedestrians for detection.

With the release of several multispectral datasets [18, 6, 45, 39], multimodal detectors have seen increasing popularity. To exploit the complementary potential between thermal and RGB images, Liu *et al.* [26] introduce a fusion method based on the Faster R-CNN. Li *et al.* [23] introduce Illumination Aware Faster R-CNN which adaptively integrates color and thermal sub-networks, and fuses the results using a weighting scheme depending on the illumination condition. Region Re-construction Network is introduced in [41] which models the relation between RGB and thermal data using a CNN. These features are then fed into a Multi-Scale Detection Network, for robust pedestrian detection. In our approach however, we use solely the thermal images and not their color counterparts.

Saliency detection. Salient object detection aims to highlight the most conspicuous object in an image and a substantial number of methods have been developed for it over the past few decades. One of the earliest works on saliency detection was presented in [21], inspired by the visual system of primates which shift focus to most conspicuous objects across the visual scene. Traditional saliency detection methods involved using methods like global contrast [5], local contrast [20] and other hand crafted features like colour and texture [29, 42]. Methods described in [17, 30] form the basis for our experiments using static saliency. A complete survey of these methods is available in [1].

Recent works typically use CNNs for salient object detection. DHSNet [27] first learns global saliency cues such as global contrast, objectness, and compactness, and then uses a novel hierarchical convolutional neural network to refine the details of the saliency maps using local context information. The use of short connections to the skip layer structure of a Holistically-Nested Edge Detector is introduced in [16]. Amulet [46] integrates multi-level features at multiple resolutions and learns to predict saliency maps by combining the features at each resolution in a recursive manner. In our experiments with deep saliency techniques, we use two state-of-the-art networks, PiCA-Net [28] and R^3 -Net [7] (explained in Section 3.2.2), to generate saliency maps from thermal images and to benchmark our salient pedestrian detection dataset.

3. Approach

In this section, we explain the task of pedestrian detection in thermal images using Faster R-CNN [32]. We then present our proposed method of augmenting thermal images with their saliency maps to improve detection performance. Finally, we describe our motivation and efforts at annotating a subset of the KAIST Multispectral Pedestrian dataset [18] at the pixel level, for use by deep saliency networks.

3.1. Baseline for Pedestrian Detection in Thermal Images using Faster R-CNN

We adapt the Faster R-CNN [32] object detector for the task of pedestrian detection in thermal images. The Faster R-CNN architecture consists of a Region Proposal Network (RPN) that is used to propose regions in an image that are most likely to contain an object, and a Fast R-CNN [14] network that classifies the objects present in that region along with refining their bounding box coordinates. Both these networks operate on shared convolutional feature maps generated by passing the input image through a backbone network (typically VGG16 [35] or ResNet101 [15]). We train the Faster R-CNN end-to-end on the thermal images from the KAIST Multispectral Pedestrian dataset [18] and present the results in Table 2.

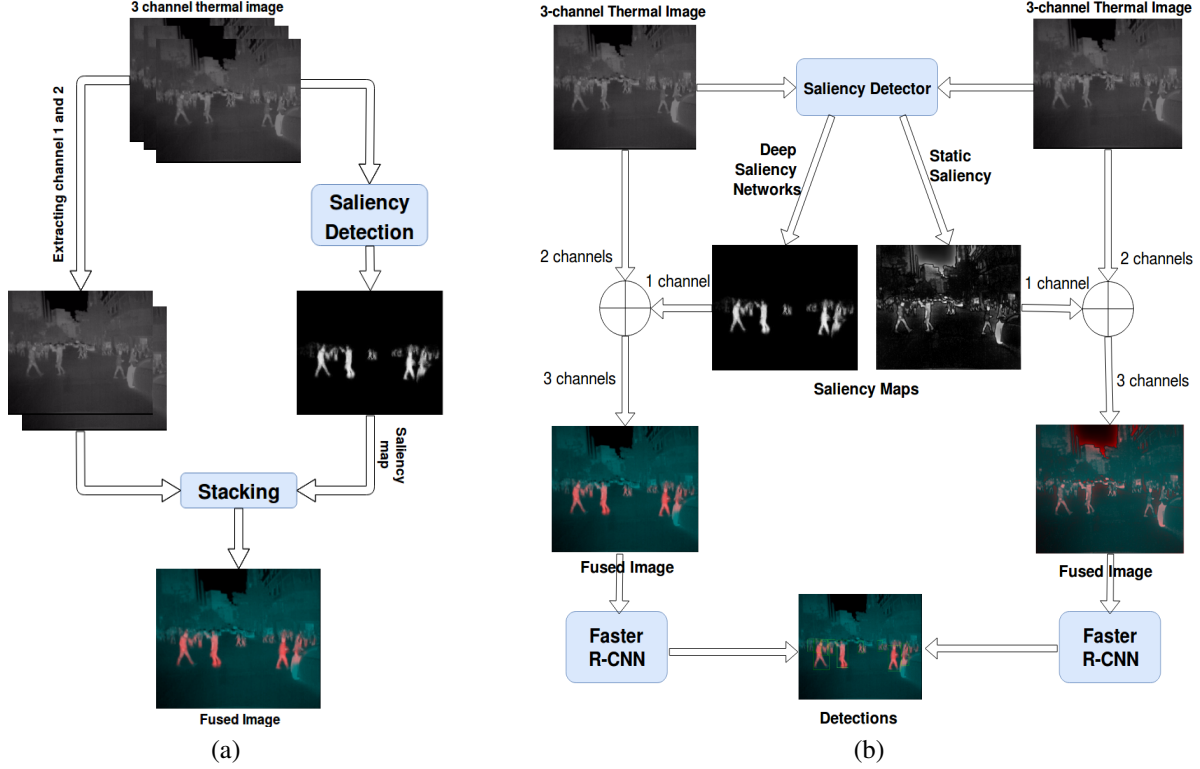


Figure 1. (a) Procedure for augmenting thermal images with saliency maps, (b) Faster R-CNN training procedure on augmented images

3.2. Our Approach: Using Saliency Maps for Improving Pedestrian Detection

We propose to use saliency maps extracted from thermal images in order to teach the pedestrian detector to “see” better through pixel level context. We expect that such a system would perform better especially during daytime when humans are more indiscernible from their surroundings in thermal images. However, saliency maps discard all textural information available in thermal images. In order to mitigate this, we augment the thermal images with their saliency maps. We do this by replacing one duplicate channel of the 3-channel thermal images with the corresponding saliency maps as shown in Figure 1(a). As seen in Figure 2, the combination of saliency maps with thermal images help illuminate the salient parts of the image, while retaining the textural information in the image. As shown in Figure 1(b), we then proceed to train the Faster R-CNN described in Section 3.1 on (i) saliency maps extracted from thermal images and (ii) thermal images augmented with their saliency maps generated using the two approaches described below.

3.2.1 Static Saliency

In this paper, we generate static saliency maps using OpenCV library [2] that uses methods described in [17] and [30]. However, the saliency maps generated using this naïve

method highlight not only pedestrians but also other salient objects in the image (as seen in Figure 2 (b) & (c)). This leaves room for a more powerful saliency detection technique that would highlight only the salient pedestrians and not any other salient objects in the image.

3.2.2 Deep Saliency Networks

We investigate two state-of-the-art deep saliency networks in this paper.

PiCA-Net [28] is a pixel-wise contextual attention network which generates an attention map for each pixel corresponding to its relevance at each location. It uses a Bidirectional LSTM to scan the image horizontally and vertically around a pixel to obtain its global context. For the local context, the attention operation is performed on a local neighboring region using convolutional layers. Finally a U-Net architecture is used to integrate the PiCA-Nets hierarchically for salient object detection.

R³ Net [7] uses a Residual Refinement Block (RRB) to learn the residuals between the ground truth and the saliency map in a recursive manner. The RRB alternatively utilizes low-level features and high-level features to refine the saliency maps at each recurrent step by adding the previous saliency map to the learned residual.

As seen in Figure 2 (d) & (e), these techniques illuminate only the pedestrians in a scene.

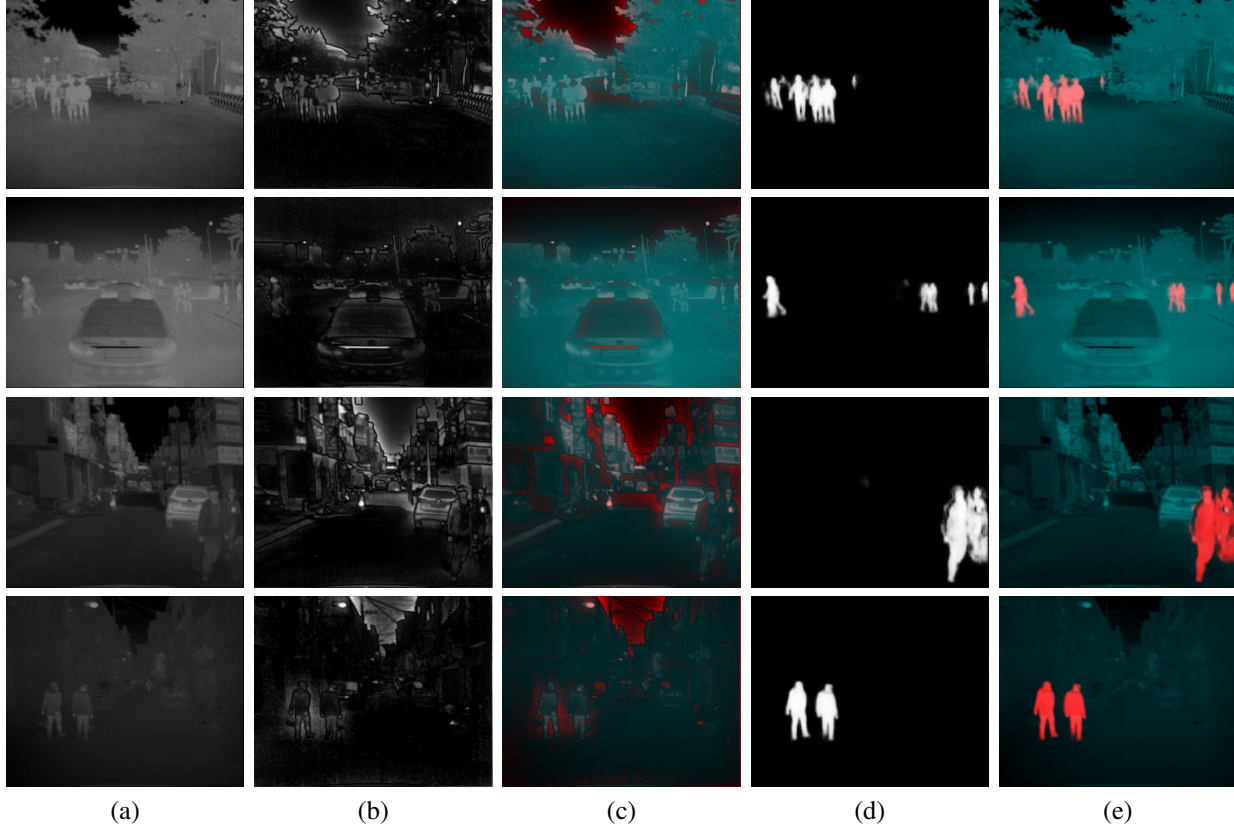


Figure 2. **Thermal images and generated saliency maps** for day (top 2 rows) and night (bottom 2 rows) images from the test set. (a) Original thermal images, (b) Static saliency maps, (c) Thermal images fused with static saliency maps, (d) Deep saliency maps, (e) Thermal images fused with deep saliency maps

3.3. Our Dataset: Annotating KAIST Multispectral Pedestrian for Salient Pedestrian Detection

In order to train a deep saliency network, we need pixel level annotations for salient objects. Since there are no publicly available thermal datasets with ground truth saliency masks for pedestrians, we create a pedestrian saliency dataset and make it publicly available¹ to facilitate further research on the use of saliency techniques for multispectral pedestrian detection.

We select 1702 images from the training set of the KAIST Multispectral Pedestrian dataset [18], by sampling every 15th image from all images captured during day and every 10th image from all images captured during night, containing pedestrians. These images were selected in order to have roughly the same number of images captured at both times of the day (913 day images and 789 night images), containing 4170 instances of pedestrians. We manually annotate these images using the VGG Image Annotator [12] tool to generate the ground truth saliency masks based on the location of the bounding boxes on pedestrians in the

original dataset. Additionally, we create a set of 362 images with similar annotations from the test set to validate our deep saliency detection networks, with 193 day images and 169 night images, containing 1029 instances of pedestrians. Figure 3 shows sample images and annotations from the new KAIST Pedestrian Saliency Detection Dataset. The distribution of pedestrians per frame in the training and test sets are shown in Figure 4. Note however that the pixel level annotations are not completely precise, so these annotations might not be suitable for fine semantic segmentation tasks. However, benchmark results in Table 1 show that this dataset works reasonably well for salient pedestrian detection tasks.

4. Experiments

4.1. Datasets and Evaluation Protocols

For training the pedestrian detectors, we use the thermal images from the KAIST Multispectral Pedestrian Dataset [18] that contains approximately 50k training images and 45k test images from videos captured during different times of the day using thermal and RGB cameras. Following the evaluation protocol in [25, 26], we sample images every

¹<https://information-fusion-lab-umass.github.io/Salient-Pedestrian-Detection/>



Figure 3. Sample annotations from our KAIST Pedestrian Saliency Dataset. Top: Original images, Bottom: Pixel level annotations

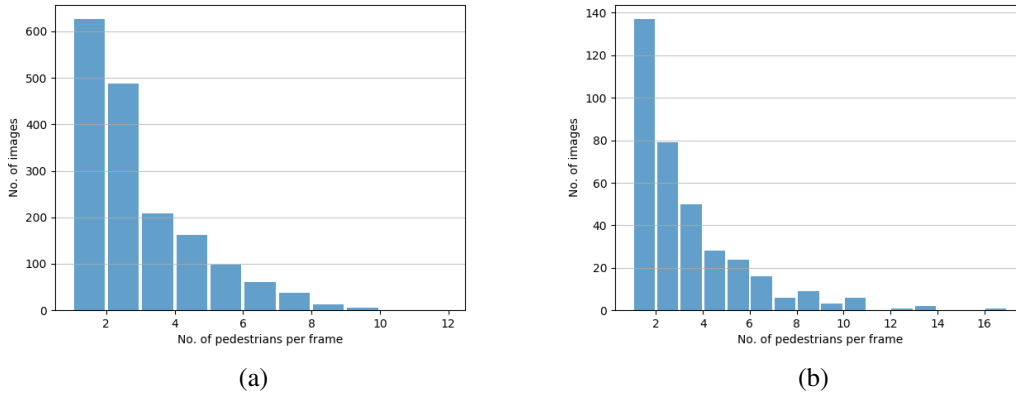


Figure 4. Distribution of Pedestrians in (a) training images (b) test images

3 frames from training videos and every 20 frames from test videos, and exclude occluded, truncated, and small (< 50 pixels) pedestrian instances. This gives us 7,601 training images (4,755 day, 2,846 night) and 2,252 test images (1,455 day, 797 night). We use the improved annotations for these 2,252 test images given in [26]. For training deep saliency networks, we annotate a subset of the KAIST Multispectral Pedestrian dataset as described in Section 3.3. Once the deep saliency networks are trained, we use them to generate saliency maps for the 7,601 training and 2,252 test images and these are then used to augment the thermal images as described in Section 3.2.

For evaluating pedestrian detection, we report the Log Average Miss Rate (LAMR) over the range $[10^{-2}, 10^0]$ against the False Positives Per Image (FPPI) under reasonable conditions [8] for day and night images. We also report the mean Average Precision (mAP) of detections at IOU=0.5 with the ground truth box. For evaluation of saliency detection, we use two metrics - F-measure score (F_β) which is a weighted harmonic mean of the precision and recall, and Mean Absolute Error (MAE) which com-

putes the average absolute per pixel difference between predicted saliency maps and corresponding ground truth saliency maps [16].

4.2. Implementation Details

4.2.1 Faster R-CNN for Pedestrian Detection

We use an open source implementation [43] of the original Faster R-CNN network with a few modifications. First, we remove the fifth max-pooling layer of the VGG16 backbone network. The original Faster R-CNN used 3 scales and 3 ratios for the reference anchors. We use 9 scales for the reference anchors, between 0.05 and 4. The Faster R-CNN network is initialized with VGG16 weights pretrained on ImageNet[33] and fine-tuned on data sources described in Section 3.2 for 6 epochs. We fix the first two convolutional layers of the VGG16 model and fine-tune the rest using SGD with momentum of 0.9, learning rate of 0.001, batch size of 1, and train our model using two NVIDIA Titan X GPUs with 12GB memory each.

4.2.2 Deep Saliency Networks

We train PiCA-Net [28] and R^3 -Net [7] on thermal images with pixel level annotations. For PiCA-Net, we use an open source implementation [19] and keep the same network architecture as described in the original paper. For training, we augment the training images with random mirror-flipping and random crops. The decoder is trained from scratch with a learning rate of 0.01 and encoder is fine-tuned with a learning rate of 0.001 for 16 epochs and decayed by 0.1 for another 16 epochs. We used SGD optimizer with momentum 0.9 and weight decay 0.0005. The entire setup is trained with a batch size of 4 on a single NVIDIA GTX 1080ti GPU. Also, since the generated saliency maps are of size 224×224 , we resize it to the original image size using Lanczos interpolation [36]. For R^3 -net we use the authors' implementation. As described, we initialize the parameters of the feature extraction network using weights from the ResNeXt [40] network. We use SGD with learning rate 0.001, momentum 0.9, weight decay 0.0005 and train for 9000 iterations using batch size of 10 on two NVIDIA Titan X GPUs with 12GB memory each.

4.3. Results and Analysis

4.3.1 Performance of Deep Saliency Networks on our KAIST Salient Pedestrian Detection dataset

We evaluate the performance of the PiCA-Net and R^3 -Net on the test set of our annotated KAIST Salient Pedestrian Detection dataset to provide a benchmark. The results are summarized in Table 1 and show reasonable saliency detection performance. Saliency masks generated using these networks can be seen in Figure 2 (d) & (e). Note that the saliency maps generated from the R^3 -Net have been post-processed using a fully-connected CRF [22] to improve coherence, resulting in the slightly better results as compared to PiCA-Net.

Method	F_β score	MAE
PiCA-Net	0.5942	0.0062
R^3 -Net	0.6417	0.0049

Table 1. Performance of deep saliency networks on our annotated test set

4.3.2 Quantitative analysis of Pedestrian Detection in Thermal Images using Saliency Maps

After evaluating the pedestrian detectors trained separately on thermal images, saliency maps, and thermal images augmented with saliency maps from different techniques, we find that the saliency maps indeed contribute to improved performance. The detector performance for each technique is summarized in Table 2, and Miss Rate vs FPPI plots are

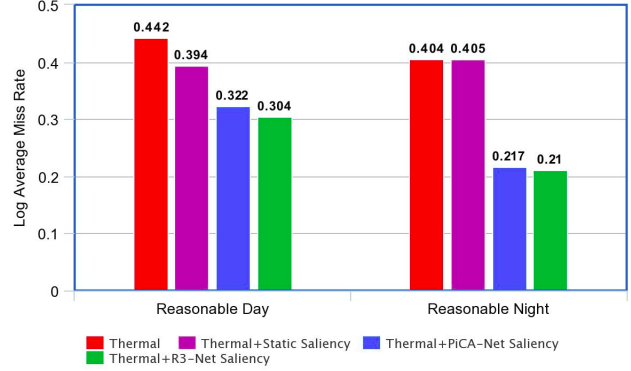


Figure 5. Comparison of Miss Rates from different models

shown in Figure 6. Below, we discuss some of the important results.

Using only Thermal Images. Our baseline detector using only thermal images achieves a miss rate of 44.2% on the day images and 40.4% on the night images displaying a large scope for improvement. It is evident from the results however, that thermal images give better performance at nighttime compared to daytime due to low contrast heat maps during the day, as seen in Figure 7(a).

Using Thermal Images with Static Saliency Maps. The pedestrian detector achieves a miss rate of 39.4% on day thermal images combined with their static saliency maps, which is an absolute improvement from the baseline by 4.8%. However, we do not notice any improvement at nighttime, and find this method to have induced a significant number of false positives hurting the precision. This indicates that although static saliency methods show some potential, they are not viable for deployment in round-the-clock applications.

Using Thermal Images with Saliency Maps generated from Deep Networks. Our approach augmenting thermal images with deep saliency maps extracted using PiCA-Net achieves a miss rate of 32.2% for day images and 21.7% for night images, which is a considerable improvement of 12% and 18.7% respectively over the baseline. The approach augmenting saliency maps from R^3 -Net achieves a miss rate of 30.4% for day images and 21% for night images, which is an even better improvement of 13.4% and 19.4% over the baseline respectively, as illustrated in Figure 5. These improvements can be explained by the visualizations in Figure 7 which shows that these methods illuminate only pedestrians in the scenes, helping the detector identify pedestrians even under difficult lighting conditions. Moreover, R^3 -Net achieves a mean Average Precision of 68.5% during daytime which is a 6.9% improvement, and 73.2% during nighttime which is a 7.7% improvement over the baseline. This suggests that deep saliency methods are useful at all times.

Testing Condition	Metric	Dataset Used						
		Thermal	Static Saliency Maps	Static Saliency + Thermal	PiCA-Net Saliency Maps	PiCA-Net Saliency + Thermal	R ³ -Net Saliency Maps	R ³ -Net Saliency + Thermal
Day	mAP	0.616	0.590	0.645	0.571	0.640	0.576	0.685
	LAMR	0.442	0.479	0.394	0.342	0.322	0.352	0.304
Night	mAP	0.655	0.605	0.641	0.639	0.676	0.585	0.732
	LAMR	0.404	0.462	0.405	0.285	0.217	0.320	0.210

Table 2. **Comparison of results from different techniques.** Our deep saliency map fused thermal images surpass all approaches in mean Average Precision (mAP) and Log Average Miss Rate (LAMR). Top 2 results are in bold.

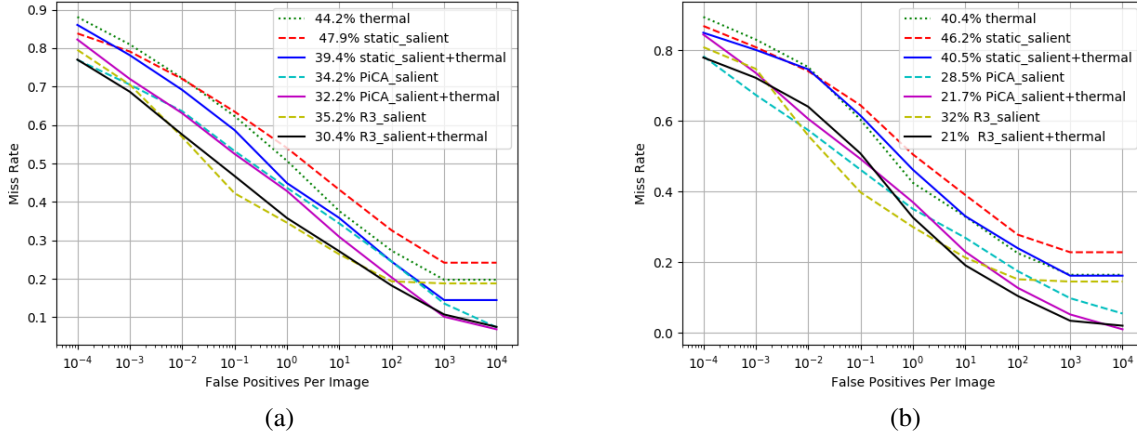


Figure 6. **Miss Rate vs FPPI curves** for a) Day reasonable conditions b) Night reasonable conditions. Our deep saliency + thermal methods are the lower curves indicating better performance compared to baseline approaches.

4.3.3 Qualitative analysis and effectiveness of saliency maps for Pedestrian Detection

Figure 7 shows detections on 4 images in different settings using all techniques. In image 1, we can see that augmenting saliency map 1(b) helps capture the rightmost missed detection in 1(a), showing its potential in cluttered scenes. In image 2(a), we see a tree detected as a false positive in the thermal image, which is a frequently occurring phenomenon in our observations. Note that the saliency maps in 2(d) & (f) puts very little emphasis on this region. Therefore, after combining the thermal image with the saliency map, the detector is able to get rid of this false positive (see 2(c), (e) & (g)). Image (3) shows comparable performance of thermal and saliency detection methods at night-time. Note that the center-right detection missed in the saliency map in 3(d) was captured in 3(e) after including the thermal information. In Image 4, the car tail-light captured by the saliency map in 4(d) is removed with the help of information from the thermal image in 4(e); whereas the detection in the middle missed by 4(a) is captured in the deep saliency maps in 4(d) & (f) and therefore included in the final detections in 4(e) & (g). This emphasizes the complementary nature of the two techniques, thus confirming our hypothesis.

5. Conclusion and Future Work

We make two important contributions in this paper. First, we provide pixel level annotations of pedestrian instances on a subset of the KAIST Multispectral Pedestrian dataset. Second, we show that deep saliency networks trained on this dataset can be used to extract saliency maps from thermal images, which when augmented with thermal images, provide complementary information to the pedestrian detector resulting in a significant improvement in performance over the baseline approach.

In this paper, we augmented thermal images with their saliency maps through a channel replacement strategy prior to feeding them into the network. It would be interesting to see if infusing the saliency map into shared layers in the network using a saliency proposal stage, and then jointly learning the pedestrian detection and the saliency detection task similar to SDS R-CNN[3] would improve the detector performance. Deep saliency techniques would also benefit from the presence of large amounts of pixel level annotations, indicating a necessary expansion of our dataset. Moreover, saliency techniques used for thermal images are also expected to work for color images and our annotations can be used for the same purpose.

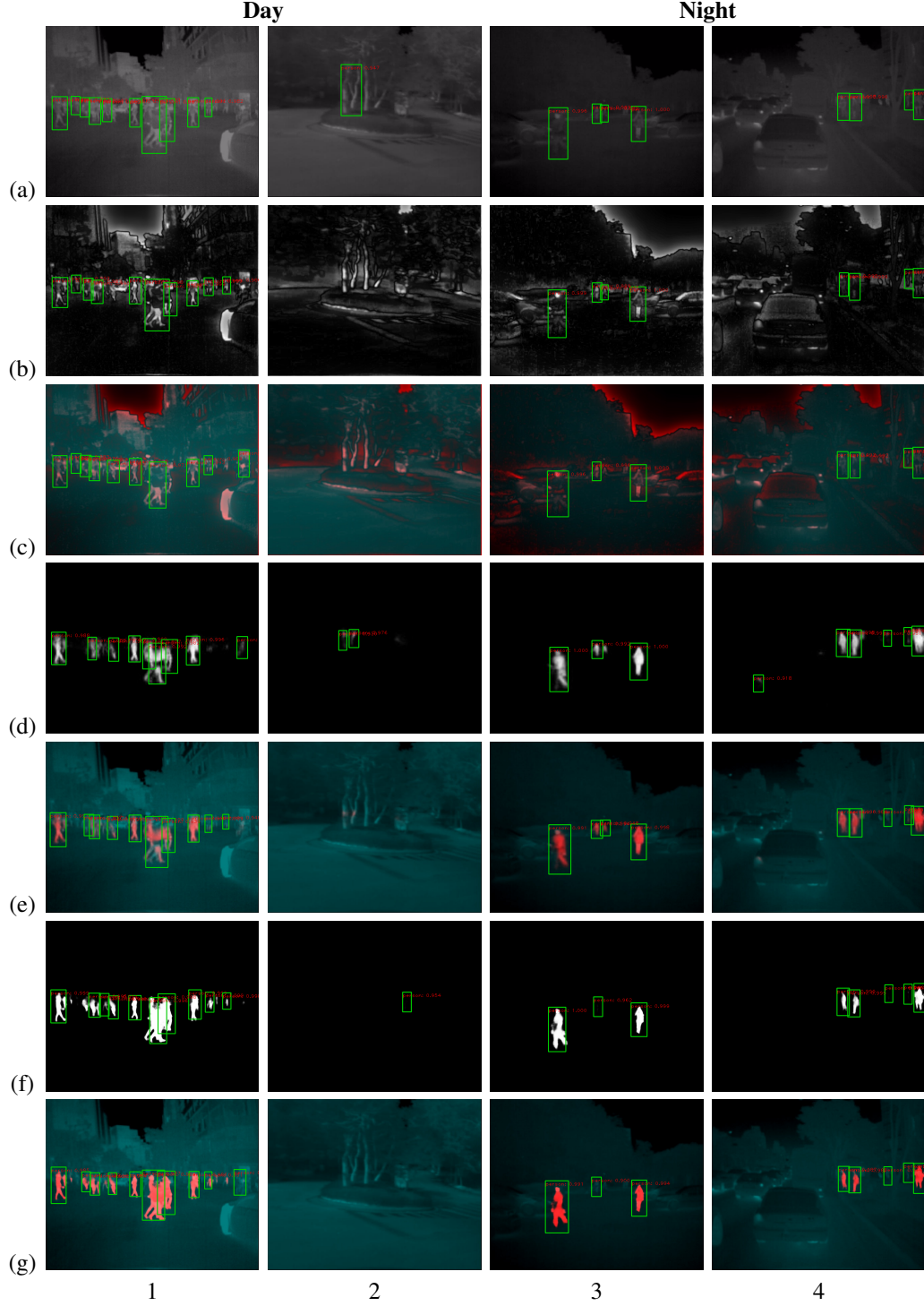


Figure 7. **Sample results from pedestrian detection** on images (1)-(4) from methods: (a) Thermal Images, (b) Static Saliency, (c) Static Saliency + Thermal, (d) PiCA-Net Saliency, (e) PiCA-Net Saliency + Thermal, (f) R^3 -Net Saliency, (g) R^3 -Net Saliency + Thermal

Acknowledgements

We would like to thank our peers who helped us improve our paper with their valuable inputs and feedback, in no

particular order - Huaizu Jiang, Takeshi Takahashi, Sarim Ahmed, Sreenivas Venkobarao, Elita Lobo, Bhanu Pratap Singh, Joie Wu, Yi Fung, Ziqiang Guan, Aruni Roy Chowdhury and Akanksha Atrey.

References

- [1] A. Borji, M. Cheng, H. Jiang, and J. Li. Salient object detection: A benchmark. *CoRR*, abs/1501.02741, 2015. 2
- [2] G. Bradski. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, 2000. 3
- [3] G. Brazil, X. Yin, and X. Liu. Illuminating pedestrians via simultaneous detection & segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4950–4959, 2017. 1, 2, 7
- [4] Z. Cai, M. Saberian, and N. Vasconcelos. Learning complexity-aware cascades for deep pedestrian detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3361–3369, 2015. 1
- [5] M. Cheng, N. J. Mitra, X. Huang, P. H. S. Torr, and S. Hu. Global contrast based salient region detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(3):569–582, March 2015. 2
- [6] 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 I*, volume 1, pages 364–369. IEEE, 2005. 2
- [7] 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 *IJCAI*, 2018. 1, 2, 3, 6
- [8] 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*, 34(4):743–761, 2012. 5
- [9] P. Dollr, R. Appel, S. Belongie, and P. Perona. Fast feature pyramids for object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(8):1532–1545, Aug 2014. 2
- [10] P. Dollr, Z. Tu, P. Perona, and S. Belongie. Integral channel features. 01 2009. 2
- [11] X. Du, M. El-Khamy, J. Lee, and L. Davis. Fused dnn: A deep neural network fusion approach to fast and robust pedestrian detection. In *2017 IEEE winter conference on applications of computer vision (WACV)*, pages 953–961. IEEE, 2017. 1
- [12] A. Dutta, A. Gupta, and A. Zissermann. VGG image annotator (VIA). <http://www.robots.ox.ac.uk/vgg/software/via/>, 2016. Version: 1.0.6, Accessed:03-01-2019. 4
- [13] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? In *Proc. CVPR*, pages 3354–3361. 1
- [14] R. Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015. 2
- [15] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 2
- [16] Q. Hou, M.-M. Cheng, X. Hu, A. Borji, Z. Tu, and P. H. Torr. Deeply supervised salient object detection with short connections. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3203–3212, 2017. 2, 5
- [17] X. Hou and L. Zhang. Saliency detection: A spectral residual approach. In *In IEEE Conference on Computer Vision and Pattern Recognition (CVPR07)*. IEEE Computer Society, pages 1–8, 2007. 2, 3
- [18] 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*, pages 1037–1045, 2015. 1, 2, 4
- [19] Y. Jaehoon. Pytorch implementation of picanet: Learning pixel-wise contextual attention for saliency detection. URL <https://github.com/Ugness/PiCANet-Implementation>, 2018. 6
- [20] D. A. Klein and S. Frintrop. Center-surround divergence of feature statistics for salient object detection. *2011 International Conference on Computer Vision*, pages 2214–2219, 2011. 2
- [21] C. Koch and S. Ullman. Shifts in selective visual attention: towards the underlying neural circuitry. In *Matters of intelligence*, pages 115–141. Springer, 1987. 1, 2
- [22] P. Krähenbühl and V. Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. In *Advances in neural information processing systems*, pages 109–117, 2011. 6
- [23] C. Li, D. Song, R. Tong, and M. Tang. Illumination-aware faster r-cnn for robust multispectral pedestrian detection. *Pattern Recognition*, 85:161–171, 2019. 1, 2
- [24] J. Li, X. Liang, S. Shen, T. Xu, J. Feng, and S. Yan. Scale-aware fast r-cnn for pedestrian detection. *IEEE Transactions on Multimedia*, 20(4):985–996, 2018. 1, 2
- [25] J. Liu. *Exploiting multispectral and contextual information to improve human detection*. Rutgers The State University of New Jersey-New Brunswick, 2017. 4
- [26] J. Liu, S. Zhang, S. Wang, and D. N. Metaxas. Multispectral deep neural networks for pedestrian detection. *arXiv preprint arXiv:1611.02644*, 2016. 1, 2, 4, 5
- [27] N. Liu and J. Han. Dhsnet: Deep hierarchical saliency network for salient object detection. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 678–686, June 2016. 2
- [28] N. Liu and J. Han. Picanet: Learning pixel-wise contextual attention in convnets and its application in saliency detection. *CoRR*, abs/1708.06433, 2017. 1, 2, 3, 6
- [29] T. Liu, Z. Yuan, J. Sun, J. Wang, N. Zheng, X. Tang, and H. Shum. Learning to detect a salient object. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(2):353–367, Feb 2011. 2
- [30] S. Montabone and A. Soto. Human detection using a mobile platform and novel features derived from a visual saliency mechanism. *Image and Vision Computing*, 28(3):391–402, 2010. 2, 3
- [31] W. Nam, P. Dollár, and J. H. Han. Local decorrelation for improved pedestrian detection. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 1, NIPS'14*, pages 424–432, Cambridge, MA, USA, 2014. MIT Press. 2

- [32] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. 1, 2
- [33] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. S. Bernstein, A. C. Berg, and F. Li. Imagenet large scale visual recognition challenge. *CoRR*, abs/1409.0575, 2014. 5
- [34] P. Sermanet, K. Kavukcuoglu, S. Chintala, and Y. LeCun. Pedestrian detection with unsupervised multi-stage feature learning. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 3626–3633, 2013. 2
- [35] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. 2
- [36] K. Turkowski. Filters for common resampling tasks. In *Graphics gems*, pages 147–165. Academic Press Professional, Inc., 1990. 6
- [37] J. Wagner, V. Fischer, M. Herman, and S. Behnke. Multi-spectral pedestrian detection using deep fusion convolutional neural networks. In *24th European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning (ESANN)*, pages 509–514, 2016. 1
- [38] X. Wang, M. Wang, and W. Li. Scene-specific pedestrian detection for static video surveillance. *IEEE transactions on pattern analysis and machine intelligence*, 36(2):361–374, 2014. 1
- [39] Z. Wu, N. Fuller, D. Theriault, and M. Betke. A thermal infrared video benchmark for visual analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 201–208, 2014. 2
- [40] S. Xie, R. Girshick, P. Dollr, Z. Tu, and K. He. Aggregated residual transformations for deep neural networks. *arXiv preprint arXiv:1611.05431*, 2016. 6
- [41] 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*, pages 5363–5371, 2017. 1, 2
- [42] C. Yang, L. Zhang, H. Lu, X. Ruan, and M. Yang. Saliency detection via graph-based manifold ranking. In *2013 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3166–3173, June 2013. 2
- [43] J. Yang, J. Lu, D. Batra, and D. Parikh. A faster pytorch implementation of faster r-cnn. <https://github.com/jwyang/faster-rcnn.pytorch>, 2017. 5
- [44] L. Zhang, L. Lin, X. Liang, and K. He. Is faster r-cnn doing well for pedestrian detection? 07 2016. 2
- [45] M. M. Zhang, J. Choi, K. Daniilidis, M. T. Wolf, and C. Kanan. Vais: A dataset for recognizing maritime imagery in the visible and infrared spectrums. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 10–16, 2015. 2
- [46] P. Zhang, D. Wang, H. Lu, H. Wang, and X. Ruan. Amulet: Aggregating multi-level convolutional features for salient object detection. pages 202–211, 10 2017. 2