# Enabling energy efficient machine learning on a Ultra-Low-Power vision sensor for IoT

Francesco Paissan[1], Massimo Gottardi[2], Elisabetta Farella[1]

[1]*E3DA - ICT-irst,* [2]*IRIS - CMM*

*Fondazione Bruno Kessler*

Trento, Italy

{fpaissan, gottardi, efarella}@fbk.eu

*Abstract*—The Internet of Things (IoT) and smart city paradigm includes ubiquitous technology to extract context information in order to return useful services to users and citizens. An essential role in this scenario is often played by computer vision applications, requiring the acquisition of images from specific devices. The need for high-end cameras often penalizes this process since they are power-hungry and ask for high computational resources to be processed. Thus, the availability of novel low-power vision sensors, implementing advanced features like in-hardware motion detection, is crucial for computer vision in the IoT domain. Unfortunately, to be highly energy-efficient, these sensors might worsen the perception performance (e.g., resolution, frame rate, color). Therefore, domain-specific pipelines are usually delivered in order to exploit the full potential of these cameras. This paper presents the development, analysis, and embedded implementation of a real-time detection, classification and tracking pipeline able to exploit the full potential of background filtering Smart Vision Sensors (SVS). The power consumption obtained for the inference - which requires 8ms - is 7.5 mW.

## I. INTRODUCTION

Vision Sensor Nodes (VSN) are the building blocks of Wireless Vision Sensors Networks (WVSN), a promising technology for its extended application domain that ranges from structural health analysis [1] to urban scenarios observation [2] [3]. While developing an algorithm for a VSN, it is crucial to consider the amount of data transmitted, which can be prohibitive [4] and have a high impact on energy efficiency. Literature suggests two viable solutions: bandwidth compression [5] or hardware solutions such as the near-sensor image processing paradigm [6] [7]. For the latter, it is crucial to reduce the amount of redundant information to be processed; therefore, event-based Smart Vision Sensors (SVS) enabling in-hardware motion detection algorithms are used [8] [9].

To exploit the full potential of these sensor technologies - which usually are also characterised by lower quality images (e.g. lower spatial and pixel resolution) - application-specific processing pipelines are the most promising option, offering a trade-off between performance and power consumption [10] [11] [12].

To prove the efficiency of these pipelines, we implemented a detection, classification and tracking pipeline on images from a low-power event-based VGA sensor with background subtraction and basic analysis capabilities. Therefore, the contributions of the work are: in
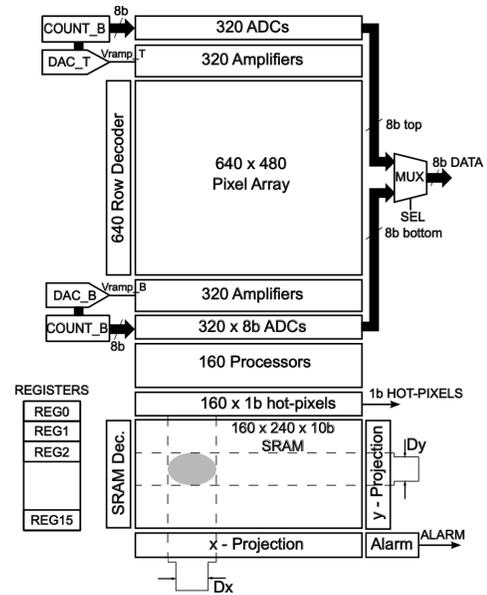


Fig. 1: Block diagram of the vision sensor generating an alert signal in presence of a moving target.

- the design of a sub-optimal detector targeting low operations count;
- the design of a pipeline targeting resource-constrained platforms for the IoT, implementing detection, classification and tracking of cars and pedestrians;
- MCU implementation and power consumption analysis of the proposed approach.

## II. RELATED WORK

Camera-based solutions for object detection already exist but, as already mentioned, are costly from a computational and energy point of view [13]. In recent years, research focused on investigating alternatives to GP-GPUs clusters to fit resource-constrained platforms for applications in robotics, smart maintenance, and Internet of Things (IoT) [14] [15]. However, many of these applications use off-the-shelf processing units that still have a power demand around $5\,\mathrm{W}$, more than what a solar-powered node can harvest. This system's high power consumption is to be attributed to the computational power needed to analyze high-resolution image streams. Smart Vision

DATE Friday Workshop on *System-level Design Methods for Deep Learning on Heterogeneous Architectures (SLOHA 2021)*, Virtual Workshop, February 5, 2021.

1

Sensors aim at reducing the amount of redundant information in images, thus enabling more efficient processing of the scene. Technically, the most popular approach to detect and remove redundant information before delivering it to the processing unit is frame difference. In [9], the authors implemented an advanced frame difference technique that exploits temporal averaging to augment the range of detectable objects. Also [16], [17] and [18] exploit the frame difference paradigm with some minor adaptations (e.g. dynamic thresholding). More advanced techniques as the one implemented in [19], which exploits a 64x64 CMOS image sensor with on-chip moving object detection and rough localization, lack flexibility since the cluster for frame-recognition cannot be changed. One main drawback of this frame difference-based SVSs is that an event might be triggered in the presence of small movements of non-static objects as tree branches or sea waves.These are not useful targets for the proposed application, thus leading to higher power consumption.

Given the many advantages in exploiting the near-sensor processing paradigm, many studies on IoT-compatible processing pipelines as developed both for binary and RGB images. Binary image processing were of high interest in the last decades [20] [21] and afterward abandoned due to the improvements in the imaging and processing field. Nonetheless, this algorithms were rediscovered for the processing of ultra-low-power smart visual sensors generated images, where the binary output enables smart strategies to reduce power consumption and computational cost. [10] present a image processing node that implements Binarized Neural Networks (BNNs) to process the output of a SVS more efficiently. In [22] the authors exploit the synergy between smart sensors and energy-efficient quad-core cluster processor for data analysis. The evaluation setup is indoor and therefore it does not exploit the potential of SVSs, like the robustness to small movements, typical of outdoor environments. Also in [11] the authors exploit BNNs to filter pedestrian-related events, thus reasoning on higher order information for triggering event alarms.

## III. VISION SENSOR

In this work, we used a low-power VGA vision sensor embedding motion detection based on dynamic background subtraction. As shown in Fig. 1, the algorithm is applied on a QQVGA sub-sampled image, generating a motion bitmap of 120 x 160 pixel, where asserted pixels are called Hot-Pixels (HP). The bitmap is de-noised through a programmable erosion filter and used to generate two projection vectors along the *x* (160 pixels) and *y* (120 pixels) axes of the array. The size and aspect ratio of the moving target to generate the alert are assigned by setting constraints on the number of contiguous motion pixels in the (x,y)-projection vectors. The vision sensor operates in Motion Detection (MD), looking for events with no off-chip data delivering, and Imaging Mode (IM), activated after the alert, delivering VGA gray-scale image together with the related motion bitmap. A low-power FPGA controls the sensor and generates the interface with the external processor. After an event is detected, the sensor sets the ALARM to
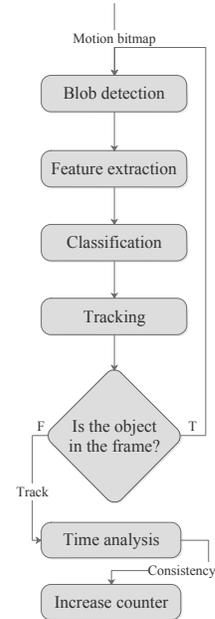


Fig. 2: Detection, classification and tracking pipeline.

trigger the processor and starts delivering 10 consecutive frames, sufficient to execute visual tasks such as detection, classification and tracking of the objects that generated the alarm. The frames are delivered as gray-level VGA, together with the sub-sampled motion bitmap encoded into the least significant bit of the VGA pixels.
When available, the frames are delivered from the FPGA to a low-power MCU through a DCMI bus.

Because the sensor is still in a prototyping stage, we are forced to sub-sample the images from the gray-scale VGA output, therefore adding overhead to the processing stage caused by the acquisition of a 4x bigger frame and the extraction of the motion information. As soon as an application-specific interface for the sensor will be ready, this operation will be avoided and the MCU would receive as independent signals the (x,y)-projections of the alarm-triggering object, therefore removing even more pre-processing from the vision pipeline.

## IV. PROCESSING PIPELINE

We implemented a pipeline for detection, classification and tracking of pedestrians and cars. For this, we exploited the features of the SVS described in the previous section. Our work is to be considered as an extension of the pipeline proposed for detection and classification in [23].

We refer to Figure 2 for an illustration of the algorithm presented in this work.

### A. Detection

Object detection - as intended in the computer vision domain - is the task of extracting instances of semantic objects (belonging to a certain class) from images. Our detection algorithm is developed for binary images, encoded as 1-pixels (active pixels) when movement is detected and 0-pixels (inactive pixels) where no movement is detected.
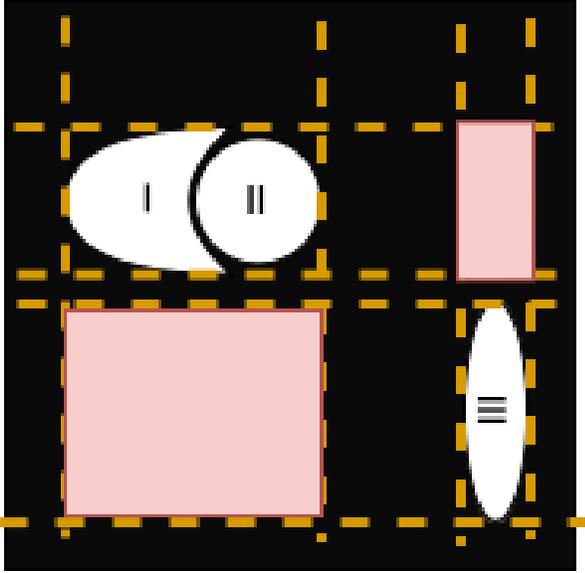
Fig. 3: Downsides on the proposed region proposal algorithm. Red rectangles are extra area which will be removed in the feature extraction process. Blob II is not visible.

As regard the detection sub-task, the contribution of this work is the development of a highly specific algorithm, which aims at being as neutral as possible from a computational cost analysis with respect to the rest of the pipeline. Our approach is compared with the method implemented in [23] and proposed by [21], which is highly accurate but has a significant computational cost – almost 80% of the required operations for the detection and classification pipeline are due to the choice of the detection algorithm.

The proposed detection algorithm is a two-stage detector. In particular, we developed a region-proposal algorithm that consists of a Cartesian product and aims at reconstructing the regions that triggered the alarm from the (x,y)-projection. From Figure 3, it is clear that this approach is limited with respect to the more sophisticated approach developed by Suzuki et al. [21] in extracting higher-level information (e.g., the hierarchical structure of contours) and less false positives. In fact, using our approach the system is more likely to detect false positives and therefore a filter for void pixel areas in needed. However, given the sensor's application domain and considering the SVS output, those limitations are not a big concern. In fact, this algorithm was tested on a sequence on 800 frames from an urban scenario recorded using the SVS. The results showed that, on average, our region proposal algorithm detects one extra object with respect to Suzuki's algorithm. Moreover, we computed the intersection-over-union of the non-empty bounding boxes with the ones detected by Suzuki's algorithm and found that, on average, it is around 0.30. This is Because the region-proposal algorithm detects bounding boxes which area is around 1.5 times compared with the more sophisticated method.

Formally, to detect regions of interest (ROI) from the (x,y)-projections, we consider the falling and rising edges of the projections and compute the Cartesian product between the rising and falling edges on the two axes obtaining the upper left and lower right corners as output sets respectively. Therefore, without considering the temporary sensor's limitations, the only operations needed to run the detection stage are memory read and write, negligible with respect to operations needed to run Suzuki's algorithm. Afterwards, in the feature extraction stage, the moments of the extracted distributions are computed and used to remove empty regions from the proposed set.

Despite the significant improvement in operation count for detection, this approach does not seem to impact classification and tracking performance on average. On the other hand, in some complicated scenarios with complex binary representation 3 this approach would inevitably perform worst than Suzuki's method. But those scenarios are not characteristic of the proposed application and, given the low-resolution of the sensor, are not likely to be detectable if present.

### B. Classification

For the classification stage, we validate the model proposed in [23]. In fact, classification is based on a linear kernel Support Vector Machine (SVM), trained on an artificially created dataset with 264 instances equally distributed between the two classes. In fact, we compared classification results using different features as for example invariant moments. Finally, we concluded that the accuracy of the model did not change much and remains nearly 98% on the artificial dataset. To add value to the proposed solution, which is already good from a computational perspective, we run tests to understand how the model behaves with respect to the features, which are:

- its area;
- the variances of the distribution associated to it on both axes, which are considered as descriptors of its shape;
- the minimum vertical coordinate of the bounding box including each blob, due to its correlation with respect to the scaling of the area.

As expected from a visual analysis of the sensor's output, we can see from Table I that the permutation importance [24] of the variance on the y axis is the most important features for the classification, together with the blob's area.

| Feature | Weigth |
|---|---|
| $\sigma_y^2$ | 0.38 ±0.05 |
| Area | 0.16 ±0.04 |
| Blob position | 0.11 ±0.03 |
| $\sigma_x^2$ | 0.06 ±0.03 |

TABLE I: Permutation importance of classification features.

To better understand how the model behaves with respect to the two most important features, we used a second order Accumulated Local Effects (ALE) [25] plot (Fig. 4). As we can see by analysing the extracted plot, the model prediction is mainly static all over the feature space, except in the lower-right corner. Meaning that the model output varies more in

that section of the feature space and therefore the predicted class changes presumably in that space.
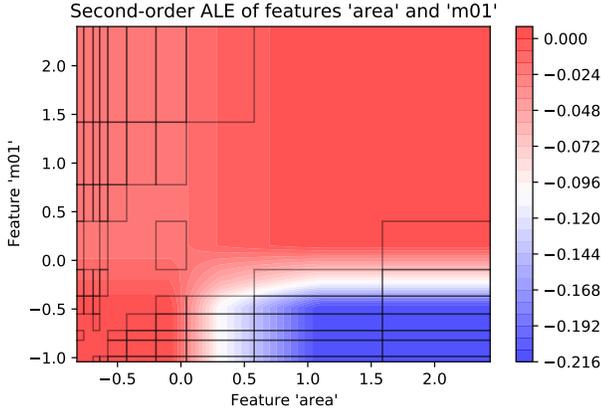


Fig. 4: Second order Accumulate Local Effects plot of features area and $\sigma_y^2$ (m01) already normalized. ALE computed on 10x10 bins.

## C. Tracking

The tracking algorithm implemented is inspired on the Simple Online Real-Time Tracking (SORT) algorithm proposed by Bewley et. al. in [26]. Particularly, the main structure is preserved and its independent blocks are modified to be more energy and computationally efficient.

*1) Detection:* This is the first stage of the tracking algorithm. In our implementation, we exploit the computational and energy efficiency of the algorithm already described.

*2) Estimation:* In this stage, a linear velocity model of the tracked object is updated during each observation of the object taken into account. This is performed using an implementation of Kalman filter [27], that exploits Single Instruction Multiple Data (SIMD). Despite we are aware that this is not the best possible representation of the object's dynamics, we observed during experiments that it works well enough to be used for our applications.

*3) Association:* This step is the core of the tracking algorithm, and consists in associating the same object in two consecutive frames. For this, many algorithms and methodologies might be involved [28], [29], but the most commonly used is the Munkres assignment algorithm (also known as the Hungarian algorithm), which proved to work well for temporary occlusions caused by moving objects.

The correlations between each blob $B_{t-1}$ in the frame $\Sigma_{t-1}$ (at time $t-1$) and each blob $B_t$ in the frame $\Sigma_t$ (at time $t$) can be represented like a bipartite graph. Specifically, the connected components in $\Sigma_{t-1}$ make up one part of the graph, and the connected components in $\Sigma_t$ make up the other part. By definition, the parts of the graph are fully connected and each edge has a weight computed as the overlap index (or Intersection over Union distance - IoU distance) between the detection boxes, including the blobs represented by the two connected nodes. Mathematically, for all $B_{t-1} \in O_{t-1}, B_t \in O_t$ the IoU

distance, and therefore the weight of the edge from $B_{t-1}$ to $B_t$ is computed as:

$$IoU(B_{t-1}, B_t) = \frac{B_{t-1} \cap B_t}{B_{t-1} \cup B_t} \qquad (1)$$

After defining the above structure, we solve for the optimal association using the Hungarian algorithm.

*4) Create or delete track:* In this stage, we insert or remove objects from the active objects list.

Every time a connected component enters the frame, if it is detected for more than $N_{hits}$ times it is considered as active and its state estimator is initialized using the last geometry recorded. Instead, if an already-tracked object is not detected for more than $T_{lost}$ frames, it is deleted from the active blob list.

*5) Time consistency and Counter:* For each tracked object, we want the algorithm to be reliable in all its stages. For this, we embed on each track a label referring to the class in which it was predicted the most of the time. This is based on the assumption that the classification might be wrong for about 10% of the samples, but should be reliable for the rest of the cases. Finally, after the object exits the frame, the consistency analysis in time is performed and the counter for the specific class is increased.

## V. RESULTS

In order to validate the presented pipeline, we performed the analysis described in this section.

### A. Validation sequences

Three sequences, for a total of 400 frames[1], were recorded at 8 fps using the development kit located at 7 meters above the ground and with a vertical orientation of $50°$. Since we wanted to validate the system at a broader level, we tried our algorithm also on a sequence taken from the PETS2009 dataset [30]. From this dataset, we selected a video in line with the aforementioned setup and we slightly adapted it in order to make it compliant with the sensor's simulator [31].

The labeling on all sequences was done manually by counting the number of cars and pedestrian appearing in the selected sequence.

### B. Tracking Performance

We use the detection and classification pipeline described, that on our sequences has an average accuracy of 85%. Specifically, for an instance to be correctly predicted in this analysis, we require that the predicted class matches the label, and that the bounding box has an IoU above 0.5.

In order to quantitatively evaluate the performance of the tracking algorithm, we computed the error for each class ($i$) as:

$$E^i = \frac{N_t^i - N_p^i}{N_t^i} \qquad (2)$$

---

[1]The authors will keep an updated version of the dataset at https://github.com/fpaissan/SVS_BS-dataset.
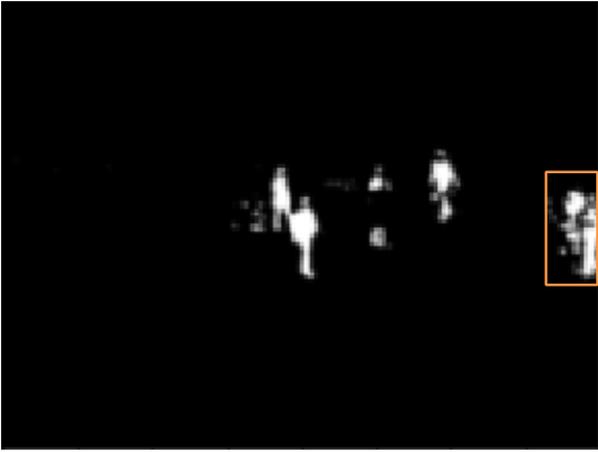
Fig. 5: Sensor output. In the orange rectangle there are two pedestrians.

where $N_t^i$ is the number of real objects belonging to class $i$ that passed in the frame, and $N_t^i$ is the number of tracked objects, always belonging to $i$. From definition, it follows that the error is zero when the tracking matches the label, negative if there are more tracked objects than labeled and positive otherwise.

TABLE II: Tracking performance over classes

| | Ground Truth | | Tracking | |
|---|---|---|---|---|
| | Car | Pedestrian | Car | Pedestrian |
| Sequence 1 | 4 | - | 4 | - |
| Sequence 2 | 1 | 2 | 1 | 2 |
| PETS2009 Seq | - | 6 | - | 5 |

As shown in Table II, the performance of the pipeline (run with $T_{lost} = 1, N_{hits} = 6$) is good, despite it has a limited set of suitable application domains. In fact, for how it is developed, it does not perform well whenever two objects are close to each other. In fact, in this case it is possible that one of the two might be considered noise - see Figure 5. Despite this, we have an average accuracy of 94.6%.

This result is interesting, since it is obtained despite the detection and classification accuracy is around 85%. Moreover, this outcome implies that by pairing the detection and classification pipeline with the tracking algorithm, we can use less accurate algorithms for detection and classification. This can help in further reducing the power consumption of the system, broadening the algorithm choice, but of course always keeping in mind that the tracking is based on the detection result. Therefore, there must be a lower bound threshold for the detection accuracy, which permits the system to work.

### C. Power consumption

To evaluate the power consumption of the proposed approach, we implemented the whole pipeline on a Cortex-M4 based board running at $80\,\mathrm{MHz}$. The power consumption of the processing, without considering the overhead caused by the sensor is of $7.5\,\mathrm{mW}$. Since the whole pre-processing due to the sensor limitation is condensed in the acquisition process, we expect a high power consumption for the acquisition stage. In fact, the measured consumption for reading 10 gray-scale frames and extract the motion bitmap is $60\,\mathrm{mW}$. However, considering smaller images with lower bit resolution (from 8-bit gray-scale to 1-bit motion bitmap) and the event-triggered nature of the system, we conclude that it can be easily powered through energy harvesting techniques (e.g. by commercial off-the-shelf solar panels).

### VI. Conclusion

In conclusion, we presented a prototype board ready to be integrated in the Internet of Things as edge vision processing unit. Targeting resource-constrained platforms, we implemented a pipeline for the detection, classification and tracking of objects, which accuracy is around 95% and real-time implementation on a microcontroller consumes $7.5\,\mathrm{mW}$ with only $8\,\mathrm{ms}$ processing window. This is achieved mainly by exploiting the SVS capabilities and adapting the algorithm to the application specific input. This is a first step in our research on bringing the computation at the very edge, while using energy efficient but low-resolution vision sensors. Many aspects of the presented approach can be extended with the idea of exploiting more sophisticated hardware and software architectures, such as Convolutional Neural Networks (CNNs) and self-attention layers. In fact, we are working on an depth-wise separable convolution-based network [32] architecture based on YOLOv2 [33] able to scale with hardware constraints, therefore exploiting the hardware-aware scaling paradigm.

### References

[1] D. Lydon, M. Lydon, J. del Rincón, S. Taylor, D. Robinson, E. O'Brien, and F. N. Catbas, "Development and field testing of a time-synchronized system for multi-point displacement calculation using low-cost wireless vision-based sensors," *IEEE Sensors Journal*, vol. 18, no. 23, 2018.

[2] R. Iqbal, T. Maniak, and C. Karyotis, "Intelligent remote monitoring of parking spaces using licensed and unlicensed wireless technologies," *IEEE Network*, vol. 33, no. 4, July 2019.

[3] A. Prati, C. Shan, and K. I.-K. Wang, "Sensors, vision and networks: From video surveillance to activity recognition and health monitoring," *Journal of Ambient Intelligence and Smart Environments*, vol. 11, 2019.

[4] G. Pekhteryev, Z. Sahinoglu, P. Orlik, and G. Bhatti, "Image transmission over ieee 802.15. 4 and zigbee networks," in *2005 IEEE International Symposium on Circuits and Systems.* IEEE, 2005.

[5] P. Chen, P. Ahammad, C. Boyer, S.-I. Huang, L. Lin, E. Lobaton, M. Meingast, S. Oh, S. Wang, P. Yan *et al.*, "Citric: A low-bandwidth wireless camera network platform," in *2008 Second ACM/IEEE International Conference on Distributed Smart Cameras.* IEEE, 2008.

[6] R. Forchheimer and A. Astrom, "Near-sensor image processing: A new paradigm," *IEEE Transactions on Image Processing*, vol. 3, no. 6, 1994.

[7] J.-E. Eklund, C. Svensson, and A. Astrom, "Vlsi implementation of a focal plane image processor-a realization of the near-sensor image processing concept," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 4, no. 3, 1996.

[8] M. Gottardi, N. Massari, and S. A. Jawed, "A 100 $\mu$ W 128 $\times$ 64 Pixels Contrast-Based Asynchronous Binary Vision Sensor for Sensor Networks Applications," *IEEE Journal of Solid-State Circuits*, vol. 44, no. 5, pp. 1582–1592, May 2009.

[9] G. Kim, M. Barangi, Z. Foo, N. Pinckney, S. Bang, D. Blaauw, and D. Sylvester, "A 467nw cmos visual motion sensor with temporal averaging and pixel aggregation," in *2013 IEEE International Solid-State Circuits Conference Digest of Technical Papers*. IEEE, 2013, pp. 480–481.

[10] M. Rusci, D. Rossi, E. Flamand, M. Gottardi, E. Farella, and L. Benini, "Always-on visual node with a hardware-software event-based binarized neural network inference engine." New York, NY, USA: Association for Computing Machinery, 2018. [Online]. Available: https://doi.org/10.1145/3203217.3204463

[11] F. C. Ojeda, A. Bisulco, D. Kepple, V. Isler, and D. D. Lee, "On-device event filtering with binary neural networks for pedestrian detection using neuromorphic vision sensors," in *2020 IEEE International Conference on Image Processing (ICIP)*, 2020, pp. 3084–3088.

[12] X. Ling, J. Sheng, O. Baiocchi, X. Liu, and M. E. Tolentino, "Identifying parking spaces detecting occupancy using vision-based iot devices," in *2017 Global Internet of Things Summit (GIoTS)*, 2017, pp. 1–6.

[13] Z. Cai, Q. Fan, R. S. Feris, and N. Vasconcelos, "A unified multi-scale deep convolutional neural network for fast object detection," in *European conference on computer vision*. Springer, 2016, pp. 354–370.

[14] E. Macias-Garcia, D. Galeana-Perez, and E. Bayro-Corrochano, "Cnn based perception system for collision avoidance in mobile robots using stereo vision," in *2020 International Joint Conference on Neural Networks (IJCNN)*, 2020, pp. 1–7.

[15] E. A. Castillo and A. Ahmadinia, "Iot-based multi-view machine vision systems," in *2019 IEEE International Conference on Big Data (Big Data)*, 2019, pp. 5206–5212.

[16] J. Choi, S. Park, J. Cho, and E. Yoon, "A 3.4-$\mu w$ object-adaptive cmos image sensor with embedded feature extraction algorithm for motion-triggered object-of-interest imaging," *IEEE Journal of Solid-State Circuits*, vol. 49, no. 1, pp. 289–300, 2013.

[17] O. Kumagai, A. Niwa, K. Hanzawa, H. Kato, S. Futami, T. Ohyama, T. Imoto, M. Nakamizo, H. Murakami, T. Nishino *et al.*, "A 1/4-inch 3.9 mpixel low-power event-driven back-illuminated stacked cmos image sensor," in *2018 IEEE International Solid-State Circuits Conference-(ISSCC)*. IEEE, 2018, pp. 86–88.

[18] K. D. Choo, L. Xu, Y. Kim, J.-H. Seol, X. Wu, D. Sylvester, and D. Blaauw, "5.2 energy-efficient low-noise cmos image sensor with capacitor array-assisted charge-injection sar adc for motion-triggered low-power iot applications," in *2019 IEEE International Solid-State Circuits Conference-(ISSCC)*. IEEE, 2019, pp. 96–98.

[19] B. Zhao, X. Zhang, S. Chen, K.-S. Low, and H. Zhuang, "A 64 x 64 cmos image sensor with on-chip moving object detection and localization," *IEEE transactions on circuits and systems for video technology*, vol. 22, no. 4, pp. 581–588, 2012.

[20] S. B. Gray, "Local properties of binary images in two dimensions," *IEEE Transactions on Computers*, vol. 100, no. 5, pp. 551–561, 1971.

[21] S. Suzuki and K. be, "Topological structural analysis of digitized binary images by border following," *Computer Vision, Graphics, and Image Processing*, vol. 30, no. 1, pp. 32 – 46, 1985.

[22] M. Rusci, D. Rossi, M. Lecca, M. Gottardi, E. Farella, and L. Benini, "An event-driven ultra-low-power smart visual sensor," *IEEE Sensors Journal*, vol. 16, no. 13, pp. 5344–5353, 2016.

[23] F. Paissan, G. Cerutti, M. Gottardi, and E. Farella, "People/car classification using an ultra-low-power smart vision sensor," in *2019 IEEE 8th International Workshop on Advances in Sensors and Interfaces*, 2019.

[24] A. Altmann, L. Toloşi, O. Sander, and T. Lengauer, "Permutation importance: a corrected feature importance measure," *Bioinformatics*, vol. 26, no. 10, pp. 1340–1347, 2010.

[25] D. W. Apley and J. Zhu, "Visualizing the effects of predictor variables in black box supervised learning models," 2019.

[26] A. Bewley, Z. Ge, L. Ott, F. Ramos, and B. Upcroft, "Simple online and realtime tracking," in *2016 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2016.

[27] R. E. Kalman, "A new approach to linear filtering and prediction problems," *Journal of basic Engineering*, vol. 82, no. 1, 1960.

[28] N. Wojke, A. Bewley, and D. Paulus, "Simple online and realtime tracking with a deep association metric," in *2017 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2017.

[29] S. Gupte, O. Masoud, R. F. Martin, and N. P. Papanikolopoulos, "Detection and classification of vehicles," *IEEE Transactions on intelligent transportation systems*, vol. 3, no. 1, 2002.

[30] J. Ferryman and A. Shahrokni, "Pets2009: Dataset and challenge," in *2009 Twelfth IEEE International Workshop on Performance Evaluation of Tracking and Surveillance*, Dec 2009.

[31] [Online]. Available: https://forensor-simulator.fbk.eu/static/index.html

[32] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.

[33] J. Redmon and A. Farhadi, "Yolo9000: Better, faster, stronger. arxiv 2016," *arXiv preprint arXiv:1612.08242*, 2016.